

# THE ANNALS *of* MATHEMATICAL STATISTICS

THE ANNALS OF MATHEMATICAL STATISTICS IS AFFILIATED  
WITH THE AMERICAN STATISTICAL ASSOCIATION AND IS  
DEVOTED TO THE THEORY AND APPLICATION OF  
MATHEMATICAL STATISTICS

EDITORIAL COMMITTEE

H. C. CARVER  
A. L. O'TOOLE  
T. E. RAIFORD

Volume VI, 1935

PUBLISHED QUARTERLY  
ANN ARBOR, MICHIGAN

*The Annals is not copyrighted: any articles or tables appearing therein may  
be reproduced in whole or in part at any time if accompanied by  
the proper reference to this publication*

*Four Dollars per annum*

*Made in United States of America*

Address: ANNALS OF MATHEMATICAL STATISTICS  
Post Office Box 171, Ann Arbor, Michigan



COMPOSED AND PRINTED AT THE  
WAVERLY PRESS, INC.  
BALTIMORE, MD.







## SOME INTERESTING FEATURES OF FREQUENCY CURVES

BY RICHMOND T. ZOCH

### Introduction

It is well known that in the normal error curve the points of inflection are equidistant from the mode. However it has never been pointed out that this is also a characteristic of all of the bell-shaped Pearson Frequency Curves. This fact can be most easily shown by placing the mode at the abscissa  $x = 0$ .

Many rough checks have been developed for use in applying the Theory of Least Squares. The second part of this paper develops a rough check on the computation for use when fitting a Pearson Frequency Curve to a set of observations. No rough checks on computation are given in textbooks on Pearson's Frequency Curves.

At present it is customary to follow a separate procedure for each Type of curve when computing the constants of a Pearson Frequency Curve. The third part of this paper shows how a single system may be followed for all Types. A single procedure is very desirable in order that the rough check of Part 2 may be quickly applied.

### Part 1. Points of Inflection

Perhaps nothing brings out the limitations of the bell-shaped Pearson Curves in a more striking manner than a discussion of their points of inflection. In dealing with frequency curves it is well known that any curve can be fitted to a given distribution and that the real problem in curve fitting is the selection of a curve. Figures 1, 2, and 3 illustrate three hypothetical histograms. All three of these histograms are bell-shaped yet none of them will be closely fitted by any of the Pearson Curves. The reasons will be pointed out presently.

The differential equation from which Pearson derived his system of frequency curves is

$$\frac{dy}{dx} = \frac{y(x - P)}{b_2x^2 + b_1x + b_0}.$$

By putting  $x - P = X$ , i.e. by placing the mode at the abscissa  $X = 0$ , this differential equation may be written:

$$\frac{dy}{dX} = \frac{yX}{\pm B_2X \pm B_1X + B_0}$$

where the  $+$  or  $-$  sign is taken according to the type of the curve. (It will be shown later that the constant term of the denominator must be less than zero.)

Since in the Type III curve  $B_2$  is 0 and in the "Normal Curve" both  $B_2$  and  $B_1$  are 0 it will be advantageous to consider the general case of

$$\frac{dy}{dX} = \frac{yX}{F(X)},$$

where  $F(X)$  is an integral rational function of the  $n^{\text{th}}$  degree, at once rather than considering special cases first.

If

$$\frac{dy}{dX} = \frac{yX}{F(X)},$$

then

$$\frac{d^2y}{dX^2} = \frac{y}{[F(X)]^2} \{X^2 + F(X) - XF'(X)\}.$$

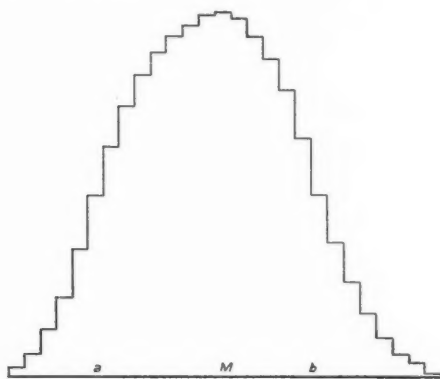


FIG. 1

In order to locate the points of inflection,  $\frac{d^2y}{dX^2}$  is equated to zero. Then we have:

$$X^2 + F(X) - XF'(X) = 0. \quad (1)$$

This equation is always of the same degree as  $F(X)$  except when  $F(X)$  is linear or constant. Hence we have proved the Theorem: If  $y = G(X)$  be the solution of the differential equation

$$\frac{dy}{dX} = \frac{yX}{F(X)},$$

then the number of points of inflection of  $y$  cannot exceed the degree of  $F(X)$  when  $F(X)$  is of degree greater than one.

Now  $F(X) = B_nX^n + B_{n-1}X^{n-1} + \dots + B_2X^2 + B_1X + B_0$ . Whence equation (1) can be written in the form:

$$(1-n)B_nX^n + (2-n)B_{n-1}X^{n-1} + (3-n)B_{n-2}X^{n-2} + \dots \\ + (r-n)B_{n-r+1}X^{n-r+1} + \dots - 3B_4X^4 - 2B_3X^3 + (1-B_2)X^2 + B_0 = 0.$$

Hence we have established the Theorem: The coefficient of the linear term of  $X$  in the equation of the points of inflection is zero.

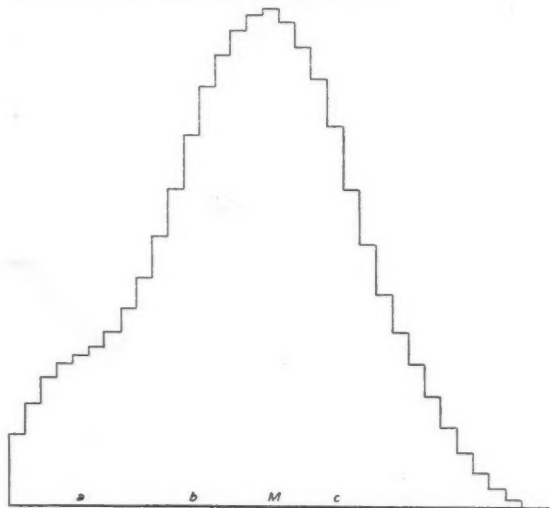


FIG. 2

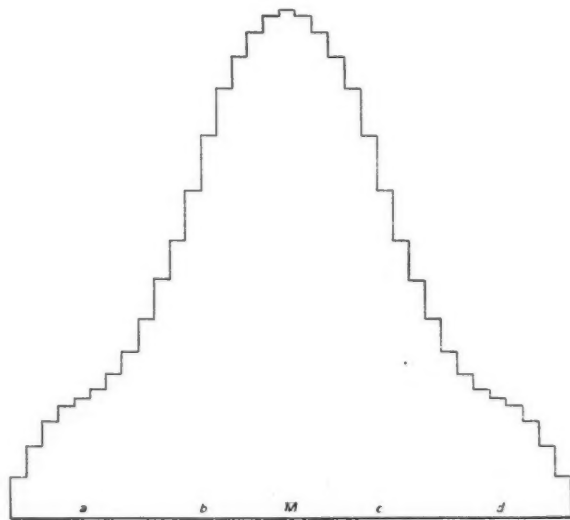


FIG. 3

For the "Normal Curve" and also for Type III,

$$B_2 = B_3 = B_4 = \dots = B_n = 0.$$

Hence the points of inflection of these two Types are given by  $X = \pm\sqrt{-B_0}$ .

For Types I and II,  $B_2$  is positive and  $B_3 = B_4 = \dots = B_n = 0$ , and the

points of inflection are  $X = \pm \sqrt{\frac{-B_0}{1 - B_2}}$ . Hence the points of inflection are undefined if  $B_2 = 1$ , are pure imaginary if  $B_2 > 1$ , and real if  $B_2 < 1$ .

For Types IV, V, VI and VII,  $B_2$  is negative and  $B_3 = \dots = B_n = 0$ , and the points of inflection are at  $X = \pm \sqrt{\frac{-B_0}{1 + |B_2|}}$ .

In some of these Types it may happen that the abscissae of the points of inflection though real will lie beyond the range of the curve. Thus Types III and VI may have 1 or 2 points of inflection, the single point of inflection occurring when  $\left| \sqrt{\frac{-B_0}{1 + B_2}} \right| >$  the range of the curve in the direction that the range is limited. Type II may have 0 or 2 points of inflection, as there will be no real points of inflection when  $B_2 \geq 1$ . Type I may have 0, 1 or 2 points of inflection. Types IV, V and VII as well as the "Normal Curve" always have 2 and only two points of inflection.

Now it should be noted that when one of the eight bell-shaped Pearson curves has two points of inflection then the abscissae of these 2 points of inflection are equidistant from the abscissa of the mode. In figure 1 a point of inflection will be at abscissa  $b$  and another at abscissa  $a$ . ( $M$  is the abscissa of the mode.) Since  $b - M \neq M - a$  none of the Pearson curves will fit this histogram closely. In figure 2, points of inflection occur at abscissae  $a$ ,  $b$ , and  $c$ . Since a Pearson curve can have at most two points of inflection no Pearson curve will fit this histogram closely. In figure 3 there are four points of inflection and no Pearson curve will fit this histogram closely.

## Part 2. Range

**DEFINITION:** A bell-shaped curve is a continuous curve which starts at zero (or zero as a limit), rises to a single maximum, at which maximum point the first derivative is zero, and then falls to zero (or zero as a limit).

Or, more formally,  $y = G(x)$  is a bell-shaped curve if  $G(x_1) = G(x_2) = 0$  and if  $G'(P) = 0$  and  $G''(P) < 0$  where  $G(x)$  is continuous and does not vanish in the interval from  $x_1$  to  $x_2$  and  $P$  is a unique point in this interval.

If a bell-shaped curve has the value of zero at two finite points, one on each side of the maximum (mode), it is said to be of limited range in both directions, or briefly, of limited range.

If a bell-shaped curve has the value of zero at only one finite point it is said to be of limited range in one direction, or also of unlimited range in one direction.

If a bell-shaped curve has the value of zero only at  $\pm \infty$ , i.e. at no finite points, it is said to be of unlimited range in both directions, or briefly, of unlimited range.

**THEOREM I:** If  $F(x)$  can be separated into a finite number of factors each either of the form  $(x - r_i)$  or  $(x^2 + 2r_i x + r_i^2 + r_{0j}^2)$  where no real root is repeated and  $y = G(x)$  is a *bell-shaped* curve which is a solution of the differential equation

$$\frac{dy}{dx} = \frac{y(x - P)}{F(x)},$$

then if  $F(x)$  has no real roots,  $y$  is of unlimited range in both directions; if all of the real roots of  $F(x)$  lie on the same side of  $P$ ,  $y$  is of limited range in one (that) direction; if at least one real root of  $F(x)$  lies on one side of  $P$  and at least one on the other side,  $y$  is of limited range in both directions.

PROOF: If  $F(x) = 0$  when  $x = P$ , we have

$$\frac{dy}{dx} = \frac{y}{g(x)}$$

where  $g(x) = F(x) \div (x - P)$ . This derivative is zero only when  $y = 0$  or  $g(x) = \pm \infty$ . Hence the solution does not have a finite maximum and therefore is not a bell-shaped curve. If  $F(x) > 0$  when  $x = P$ , we have

$$\frac{d^2y}{dx^2} = \frac{y}{[F(x)]^2} \left[ (x - P)^2 + F(x) - (x - P) \frac{d}{dx} F(x) \right]$$

and

$$\left. \frac{d^2y}{dx^2} \right|_{x=P} = \frac{y}{[F(x)]^2} [F(x)]$$

which is greater than zero and, since at a maximum the second derivative must not be greater than zero, in this case the solution would have a minimum at  $x = P$  and therefore would not be a bell-shaped curve. As the theorem concerns only those solutions which are bell-shaped curves,  $F(x) < 0$  when  $x = P$ . If  $F(x) = 0$  when  $x \neq P$  then  $\frac{dy}{dx} = \pm \infty$  unless  $y$  is also zero. Assume  $y \neq 0$ .

Since  $F(x)$  is negative, if  $y \neq 0$  when  $F(x) = 0$  then  $\frac{dy}{dx} \rightarrow -\infty$  as  $F(x) \rightarrow 0$ , for an  $x > P$ , and changes to  $+\infty$  as  $F(x)$  changes sign on passing through the value 0. Hence the curve would contain another maximum before falling to zero and therefore the solution is not a bell-shaped curve. Similar reasoning holds for an  $x < P$ . Therefore if  $y \neq 0$  when  $F(x) = 0$ , the curve is not bell-shaped. If  $y = 0$  when  $F(x) = 0$ , the curve has its range limited at this point. That is, any real number which makes  $F(x)$  vanish will also make  $y$  vanish if  $y$  represents a bell-shaped curve. Hence if all of the real roots lie on the same side of  $P$  the curve is of limited range in that direction only, while if at least one of the real roots lies on each side of  $P$  the curve is of limited range in both directions. If  $F(x)$  contains no real roots it does not vanish for any real value of  $x$ . In this case, by partial fractions the differential equation becomes:

$$\begin{aligned} \frac{dy}{y} = & \frac{k_{11} dx}{(x + r_1)^2 + r_{01}^2} + \frac{k_{21} dx}{(x + r_2)^2 + r_{02}^2} + \cdots + \frac{2k_{21}(x + r_1) dx}{(x + r_1)^2 + r_{01}^2} \\ & + \frac{2k_{22}(x + r_2) dx}{(x + r_2)^2 + r_{02}^2} + \cdots \end{aligned}$$

On integrating,

$$y = C [(x + r_1)^2 + r_{01}^2]^{k_{21}} [(x + r_2)^2 + r_{02}^2]^{k_{22}} \dots e^{k_{11} \arctan \frac{x+r_1}{r_{01}} + \dots}.$$

Hence  $y$  does not vanish for a finite real value of  $x$  and the Theorem is fully established.

**THEOREM II:** If  $F(x)$  can be separated into a finite number of factors each either of the form  $(x - r_i)$  or  $(x^2 + 2r_ix + r_i^2 + r_{0i}^2)$  where no real root is repeated and  $y = G(x)$  is a bell-shaped curve which is a solution of the differential equation  $\frac{dy}{dx} = \frac{y(x - P)}{F(x)}$ , then if  $y$  is of unlimited range,  $F(x)$  contains no real roots; if  $y$  is of limited range in one direction, all of the real roots of  $F(x)$  lie on the same (that) side of  $P$ ; if  $y$  is of limited range in both directions, at least one of the real roots of  $F(x)$  lies on one side of  $P$  and at least one on the other.

**PROOF:** By partial fractions the differential equation may be written:

$$\begin{aligned} \frac{dy}{y} = & \frac{k_{11} dx}{x - r_{11}} + \frac{k_{21} dx}{x - r_{12}} + \dots + \frac{k_{21} dx}{(x + r_{21})^2 + r_{01}^2} \\ & + \frac{k_{22} dx}{(x + r_{22})^2 + r_{02}^2} + \dots + \frac{2k_{31}(x + r_{21}) dx}{(x + r_{21})^2 + r_{01}^2} + \frac{2k_{32}(x + r_{22}) dx}{(x + r_{22})^2 + r_{02}^2} + \dots \end{aligned}$$

and on integrating:

$$y = C(x - r_{11})^{k_{11}}(x - r_{12})^{k_{12}} \dots [(x + r_{21})^2 + r_{01}^2]^{k_{31}} \dots e^{k_{21} \arctan \frac{x+r_{21}}{r_{01}} + \dots}.$$

Hence  $y = 0$  for  $x = r_{11}, r_{12}, \dots$  and for no other finite values of  $x$  provided  $k_{11}, k_{12}, \dots$  are positive. If one or more of the  $k_{ij}$  are negative,  $y = \infty$  at such points and unless some  $r_{ij}$  closer to  $P$  has previously made  $y$  vanish, the curve is not bell-shaped. Therefore, for bell-shaped curves, the exponent of the factor containing the real root of smallest absolute value on each side of  $P$  is positive. Therefore: if  $y$  is of limited range in both directions, at least one real root lies on each side of  $P$ ; if  $y$  is of unlimited range in one direction, all of the real roots lie on the same side of  $P$ ; if  $y$  is of unlimited range it contains no real roots. Hence the Theorem is established.

The effect of repeated real roots will now be considered. If a real root is repeated an odd number of times at  $x = r$ , then  $F(x)$  changes sign at  $x = r$  and the first theorem is true. If a real root is repeated an even number of times at  $x = r$ , then  $F(x)$  does not change sign at  $x = r$  and we know that either (a)  $y = 0$  at  $x = r$ ; or (b)  $y$  is finite and  $\neq 0$  and  $\frac{dy}{dx} = \pm \infty$  at  $x = r$ , i.e. there is a point of inflection at  $x = r$ . It will now be shown that (b) cannot occur. If case (b) is possible,  $y$  is continuous at  $x = r$ ,  $\frac{dy}{dx} = \pm \infty$  according as  $(r - P) \lessgtr 0$

moreover  $\frac{dy}{dx}$  does not change sign in the neighborhood of the point  $x = r$ , and  $\frac{d^2y}{dx^2}$  changes sign from  $+\infty$  to  $-\infty$  or vice versa according as  $(r - P) \lessgtr 0$ .  
Now

$$\frac{d^2y}{dx^2} = \frac{y}{[F(x)]^2} \left[ (x - P)^2 + F(x) - (x - P) \frac{d}{dx} F(x) \right].$$

Whence if  $y$  is finite and  $\neq 0$ ,  $\frac{d^2y}{dx^2}$  does not change sign at  $x = r$  because it is possible to select a neighborhood such that

$$|(x - P)^2| > \left| F(x) - (x - P) \frac{d}{dx} F(x) \right|$$

for an  $x$  differing from  $r$  by  $\epsilon$  where  $\epsilon$  is a small positive quantity. Therefore case (b) is not possible and  $y = 0$  when a real root is repeated an even number of times. That is to say the range of the curve is limited at a point where a real root is repeated an even number of times. Thus Theorem I always holds for repeated roots.

For Theorem II it is clear that this Theorem holds for repeated roots when a non-repeated root lies closer to  $P$ , and on the same side, than the repeated root. Suppose that the repeated root is the nearest root to  $P$  (on a given side of  $P$ ). Then by partial fractions:

$$\begin{aligned} \frac{dy}{y} = & \frac{k_{11} dx}{(x - r_{11})} + \frac{k_{12} dx}{(x - r_{11})^2} + \frac{k_{13} dx}{(x - r_{11})^3} + \cdots + \frac{k_{41} dx}{(x - r_{41})} + \frac{k_{42} dx}{(x - r_{42})} \\ & + \cdots + \frac{k_{21} dx}{(x + r_{21})^2 + r_{01}^2} + \frac{k_{22} dx}{(x + r_{22})^2 + r_{02}^2} + \cdots + \frac{2k_{31}(x + r_{21}) dx}{(x + r_{21})^2 + r_{01}^2} + \cdots \end{aligned}$$

and on integrating:

$$y = C(x - r_{11})^{k_{11}}(x - r_{41})^{k_{41}}(x - r_{42})^{k_{42}} \cdots [(x + r_{21})^2 + r_{01}^2]^{k_{31}} \cdots e^{k_{21} \arctan \frac{x + r_{21}}{r_{01}} + \cdots - \frac{k_{12}}{(x - r_{11})} - \frac{k_{13}}{2(x - r_{11})^2} - \cdots}$$

Hence  $y$  can = 0 only for  $x = r_{11}$  or for  $x = r_{41}, r_{42}, \cdots$  and for no other finite values of  $x$ . Since by hypothesis  $y$  is bell-shaped, then the proper  $k_{ij}$  must be positive and Theorem II always holds for repeated roots.

Theorems I and II can now be combined and generalized in the form:

**THEOREM:** If  $F(x)$  is a polynomial with real coefficients and  $y = G(x)$  is a bell-shaped curve which is a solution of the differential equation

$$\frac{dy}{dx} = \frac{y(x - P)}{F(x)},$$



then the necessary and sufficient condition: that  $y$  be of unlimited range in both directions is that  $F(x)$  have no real roots; that  $y$  be of limited range in one direction is that all of the real roots of  $F(x)$  lie on the same side of  $P$ ; that  $y$  be of limited range in both directions is that at least one real root of  $F(x)$  lie on one side of  $P$  and one on the other.

COROLLARY:  $F(x)$  must be negative throughout the range of  $y$ .

Suppose now that we have some statistics which we wish to graduate and the statistics are of such nature that we would expect a bell-shaped curve, rather than a J- or U-shaped curve, and we desire the best fit: If we use a curve which is a solution of the differential equation

$$\frac{dy}{dx} = \frac{y(x - P)}{F(x)}$$

(the Pearson Curves being special cases) to fit the statistics and if in computing the constants for the curve one of the following cases arise:

- (a)  $b'_0 < 0$  when this constant is computed,
- or (b)  $B_0 < 0$  when the origin is moved to the mode,
- or (c) a root is located within the range of the statistics then it means that:

1. A mistake may have been made in the computation; thus the Theorem just established provides a rough check on the work of computation,

2. If no mistake has been made in the computation it may indicate that the bell-shaped Pearson Curves will not closely fit the statistics and that some other graduation curves be used, e.g. the Gram-Charlier Types A or B might be tried,

3. If no mistake has been made in the computation it may happen that one of the bell-shaped Pearson Curves will give an excellent fit but a different method than or a modification of the Method of Moments should be used in order to compute the constants.

### Part 3. Computing the Constants

At present, the constants of a frequency curve are computed as follows: First the moments are computed about an arbitrary origin, then the moments about the A.M. are determined, then  $\beta_1$  and  $\beta_2$  and the criterion are computed, after which the type of curve can be selected. From this point a separate procedure is followed for each curve. Now in the above method one will not know whether a root has been located in the range of statistics or not.

Take Pearson's differential equation

$$\frac{dy}{dx} = \frac{y(x - P)}{b_2x^2 + b_1x + b_0}.$$

Put  $X = x - P$ . Then  $dX = dx$  and  $x = X + P$ , and

$$\frac{dy}{dx} = \frac{yX}{b_2(X + P)^2 + b_1(X + P) + b_0} = \frac{yX}{b_2X^2 + 2Pb_2X + b_1X + P^2b_2 + Pb_1 + b_0}.$$



Now put

$$\begin{aligned} b_2 &= B_2 \\ 2Pb_2 + b_1 &= B_1 \\ P^2b_2 + Pb_1 + b_0 &= B_0. \end{aligned}$$

Then we have

$$\frac{dy}{dX} = \frac{yX}{B_2X^2 + B_1X + B_0} \quad \text{or} \quad \frac{dy}{dx} = \frac{y(x-P)}{B_2(x-P)^2 + B_1(x-P) + B_0}. \quad (1)$$

It should be noted that for a particular curve,  $B_2$ ,  $B_1$  and  $B_0$  are constants; i.e., their values do not change with a change of the origin. The values of  $b_1$  and  $b_0$  do change with a change in the origin.

If we clear equation (1) of fractions, multiply by  $e^{yx}$  and integrate with respect to  $x$  over the range from  $x_1$  to  $x_2$ , where

$$e^{\lambda_1\eta + \frac{\lambda_2\eta^2}{2!} + \frac{\lambda_3\eta^3}{3!} + \dots} \equiv \int_{x_1}^{x_2} e^{yx} y dx,$$

then successively differentiate with respect to  $\eta$ , and equate coefficients of like powers of  $\eta$ , we finally obtain:

$$\left. \begin{aligned} \lambda_1 - P + B_1 - 2PB_2 + 2B_2\lambda_1 &= 0, \\ \lambda_2 + B_0 - PB_1 + P^2B_2 + B_1\lambda_1 - 2PB_2\lambda_1 + 3B_2\lambda_2 + B_2\lambda_1^2 &= 0, \\ \lambda_3 + 2\lambda_2B_1 - 4PB_2\lambda_2 + 4B_2\lambda_3 + 4B_2\lambda_1\lambda_2 &= 0, \\ \lambda_4 + 3B_1\lambda_3 - 6PB_2\lambda_3 + 5B_2\lambda_4 + 6B_2\lambda_2^2 + 6B_2\lambda_1\lambda_3 &= 0. \end{aligned} \right\} \quad (2)$$

Since we can compute the moments from the raw statistics and the semi-invariants from the moments, we may regard  $\lambda_2$ ,  $\lambda_3$  and  $\lambda_4$  in these equations as knowns and the  $B_0$ ,  $B_1$ ,  $B_2$ ,  $P$  and  $\lambda_1$  as unknowns. But the origin has not yet been specified. Let the origin be placed at the A.M. where  $\mu_1 = \lambda_1 = 0$ . As  $\lambda_2$ ,  $\lambda_3$ ,  $\lambda_4$ ,  $B_0$ ,  $B_1$  and  $B_2$  are unchanged by a change of origin, we have:

$$\left. \begin{aligned} B_1 - P_0 - 2P_0B_2 &= 0, \\ \lambda_2 + B_0 - P_0B_1 + P_0^2B_2 + 3B_2\lambda_2 &= 0, \\ \lambda_3 + 2B_1\lambda_2 - 4P_0B_2\lambda_2 + 4B_2\lambda_3 &= 0, \\ \lambda_4 + 3B_1\lambda_3 - 6P_0B_2\lambda_3 + 5B_2\lambda_4 + 6B_2\lambda_2^2 &= 0. \end{aligned} \right\} \quad (3)$$

Now put

$$\left. \begin{aligned} b'_0 &= B_0 - P_0B_1 + P_0^2B_2, \\ b'_1 &= B_1 - 2P_0B_2, \\ b'_2 &= B_2; \end{aligned} \right\} \quad (4)$$

then

$$\left. \begin{aligned} b'_1 - P_0 &= 0, \\ \lambda_2 + b'_0 + 3b'_2\lambda_2 &= 0, \\ \lambda_3 + 2b'_1\lambda_2 + 4b'_2\lambda_3 &= 0, \\ \lambda_4 + 3b'_1\lambda_3 + 5b'_2\lambda_4 + 6b'_2\lambda_2^2 &= 0. \end{aligned} \right\} \quad (5)$$

By reversing the transformation (4) we get:

$$\left. \begin{aligned} B_2 &= b'_2, \\ B_1 &= b'_1 + 2P_0b'_2 \\ B_0 &= b'_0 + P_0(b'_1 + P_0b'_2). \end{aligned} \right\} \quad (6)$$

Now the above theory suggests the following procedure for computing the constants of a frequency curve: First the moments are computed about an arbitrary origin, then the semi-invariants are computed (or alternatively the moments about the A.M., either step involves about the same amount of work), then the equations (5) are solved and then by means of equations (6) the  $B_2$ ,  $B_1$  and  $B_0$  are computed. Next solve the quadratic equation

$$B_2X^2 + B_1X + B_0 = 0.$$

The character of the roots of this equation indicates which type to use and it is unnecessary to compute the criterion. The constants of the frequency curve are simple functions of the roots of the above quadratic equation and can be readily found by integrating the diff. eq. (1) being careful to write the solution as a function of  $X = x - P$ . The rough checks mentioned in Part 2 can be quickly and conveniently applied when this procedure is followed.

GEORGE WASHINGTON UNIVERSITY.

## A RECONSIDERATION OF SHEPPARD'S CORRECTIONS

BY W. T. LEWIS<sup>1</sup>

In computing the moments of a frequency distribution it is customary to find first what are known as the raw moments. These are obtained on the assumption that all the material of each class interval is concentrated at the middle point of the interval. It introduces what is called a grouping error because in fact the material does not all lie at the middle point. To compensate for this error W. F. Sheppard<sup>2</sup> derived a set of corrections. The hypothesis underlying his method is that the distribution may be regarded as similar to one to which the Euler-MacLaurin summation formula without its end terms may be applied. He presupposed such a curve, found its true moments, and then the raw moments that would be obtained if its area were concentrated at several equidistant abscissae. The relationship between these raw moments and the true moments of the curve furnished him with the corrections required for that distribution. If now our observed distribution may be supposed to be sufficiently like that one, we may use his corrections also on the observed data. One may note four points of criticism.

(1) The given distribution may not be similar to the one suggested, in the sense that it would be close to such a curve if the intervals of grouping were made very small; or at all events the purpose of finding the moments may be in part to decide whether or not it would become such a curve, and so one would not like to assume that to be true at the outset. A special case of importance in which this last is true occurs when one is finding the moments of a sample in order to determine whether it may have been drawn from a presupposed universe. It is inexact to use raw moments but it is illogical to use corrections that have been proved only for the universe being tested.

(2) Sheppard's argument does not make use of the one certain fact that is given in the hypothesis, viz: that the partial area of the given distribution over each class interval is exactly as stated. In fact, if, following the argument of some authors, the given curve be assumed to be exponential, it obviously cannot have partial areas everywhere exactly equal to the several given frequencies, for in particular its partial area is not zero beyond the given range.

(3) It is common to find distributions which do not have high contact at the ends of the range and for them Sheppard's corrections certainly fail. To obviate this criticism new corrections have been derived by Pairman and Pear-

<sup>1</sup> With the assistance of Burton H. Camp.

<sup>2</sup> The true values are given on page 220 of "Mathematical Part of Elementary Statistics, by Camp, D. C. Heath and Company, 1931.

son for the so-called abrupt cases. These new corrections are adequate to care for the abrupt cases but involve so much computation that it is a fair question whether it would not be simpler, first to distribute the given material over each interval by a smoothing process, and then to find without corrections the moments of the smoothed distribution.

(4) Even if one admits Sheppard's method in general, waiving the dubious question as to whether it is proper to start with an assumed curve instead of starting with the given distribution, it is doubtful whether there are any curves which have exactly the properties required. The high contact hypothesis may be put in different language as follows: using the notation of the Handbook<sup>3</sup> page 92, let  $f(x)$  be the curve and  $x_i$  be the middle point of the slice. It is assumed that

$$\sum_i c x_i^r f^{(i)}(x_i) = \int_{-\infty}^{\infty} x^r f^{(i)}(x) dx; \quad i = 0, 1, \dots; \quad r = 0, 1, \dots;$$

$c$  being the class interval. This means that if the moments of the curve be found by using *mid-ordinates times class interval*, instead of *areas*, one will obtain exactly the true moments of the curve, and that this will remain true for all the curves which are derivatives of this curve. This property is certainly not true of the normal curve; but it is almost true when  $r$  and the class interval are both small, and it is probably due to this fact that Sheppard's corrections seem to be good in practice.

Moreover, this high contact hypothesis cannot be true for any function over a limited range if the function is developable in Taylor's series about one end of the range. For the only function which has the required properties is identically zero, since the function and all its derivatives are required to vanish at that end of the range.

The primary purpose of this paper, therefore is to derive corrections similar to Sheppard's with a different set of assumptions. The results may be used as an approximate substitute for both Sheppard's and Pairman's. That is, they will apply approximately to both extreme cases and to the intermediate cases; on the whole they give better results than Sheppard's and are not so difficult to administer as Pairman's.

The argument runs as follows. When a distribution is given merely by class intervals, there is no way of knowing exactly what the distribution would have been had the class intervals been smaller; we do not know that we have a sample from an exponential curve, and even if we did we would not know that this sample would lie close to the exponential in form. We shall, however, try to draw a graduating curve in such a manner that (a) its partial area over each class interval will equal the frequency of the given distribution over that interval; and (b) its form within each class interval will be such that it will pass smoothly into the adjacent portions to the right and left. A good way to do this is by a

<sup>3</sup> H. L. Rietz, "Handbook of Math. Stat." Houghton Mifflin Co. (1924).

freehand graph, frankly recognizing that there are many forms that will do equally well. To obtain a numerical result it is necessary to use the equation of some curve. Again frankly recognizing that there are many types which will do equally well we choose the simplest to handle:

$$y = a + bt + ct^2.$$

Let the relative frequency distribution be defined by  $f(i)$ ,  $-m \leq i \leq n$ ,  $m, n, i$  being integers. To satisfy (a) we have the equation

$$\int_{i-\frac{1}{2}}^{i+\frac{1}{2}} y dt = f(i).$$

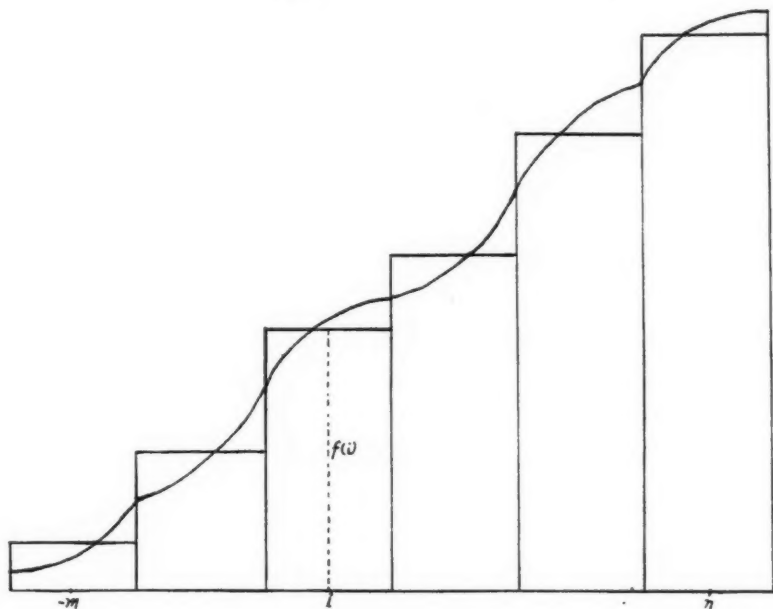


FIG. 1

To satisfy (b) we shall let

$$y = \frac{1}{2}[f(i) + f(i+1)] \text{ if } t = i + \frac{1}{2}.$$

The latter will hold for all values of  $i$  from  $-m$  to  $n-1$  inclusive, but the end intervals require special treatment. Here in order to satisfy as well as possible both the high contact and the abrupt cases, we wish to let the material be distributed according to the way the curve is behaving over the two nearest intervals on the right (at  $n$ ) or left (at  $-m$ ) rather than by the addition of zero frequencies beyond the given limits. To do this we let the slope of the parabolas be zero at the extremes:

$$\frac{dy}{dt} = 0 \quad \text{at } t = -m - \frac{1}{2} \text{ and } t = n + \frac{1}{2}.$$

Then, if for example the frequencies are increasing as one nears the right end interval, the curve will rise over the right end interval; if they are decreasing, it will fall. These three conditions are sufficient to determine a continuous curve of the sort indicated in the figure. The exact moments of the curve may be found by integration and expressed in terms of the raw moments. The details are tedious and of an elementary nature and will be given only for the mean value  $\bar{v}_1$ .

To determine the coefficients of the parabola  $y = a + bt + ct^2$  for the rectangle at  $t = i$  we may write the following three equations; the first complying with the requirement that the area under the parabola from  $t = i - \frac{1}{2}$  to  $t = i + \frac{1}{2}$  equals the area of the rectangle at  $t = i$ , the second and third giving the ordinates at  $i - \frac{1}{2}$  and  $i + \frac{1}{2}$  respectively:

$$\begin{aligned} f(i) &= \int_{i-\frac{1}{2}}^{i+\frac{1}{2}} (a + bt + ct^2) dt, \\ \frac{f(i) + f(i+1)}{2} &= a + b(i + \frac{1}{2}) + c(i + \frac{1}{2})^2, \\ \frac{f(i) + f(i-1)}{2} &= a + b(i - \frac{1}{2}) + c(i - \frac{1}{2})^2. \end{aligned}$$

Solving these three simultaneous equations we get for  $a$ ,  $b$ , and  $c$ :

$$\begin{aligned} a &= (\frac{5}{4} - 3i^2) f(i) + (\frac{3i^2}{2} - \frac{i}{2} - \frac{1}{8}) f(i+1) + (\frac{3i^2}{2} + \frac{i}{2} - \frac{1}{8}) f(i-1), \\ b &= 6if(i) + (\frac{1}{2} - 3i) f(i+1) - (\frac{1}{2} + 3i) f(i-1), \\ c &= -3f(i) + \frac{3}{2} f(i+1) + \frac{3}{2} f(i-1), \end{aligned}$$

and these hold for  $-m+1 \leq i \leq n-1$ .

For the parabola  $y = a_1 + b_1t + c_1t^2$  over the first rectangle, i.e., where  $i = -m$ , we get the equations:

$$\begin{aligned} f(-m) &= \int_{-m-\frac{1}{2}}^{-m+\frac{1}{2}} (a_1 + b_1t + c_1t^2) dt, \\ \frac{f(-m) + f(-m+1)}{2} &= a_1 + b_1(-m + \frac{1}{2}) + c_1(-m + \frac{1}{2})^2, \\ b_1 + 2c_1(-m - \frac{1}{2}) &= 0, \end{aligned}$$

and their solutions:

$$\begin{aligned} a_1 &= \frac{3}{4} (m^2 + m - \frac{1}{12}) f(-m+1) - \frac{3}{4} (m^2 + m - \frac{17}{12}) f(-m), \\ b_1 &= \frac{3}{4} (2m+1) f(-m+1) - \frac{3}{4} (2m+1) f(-m), \\ c_1 &= \frac{3}{4} f(-m+1) - \frac{3}{4} f(-m). \end{aligned}$$

Similarly for the parabola  $y = a_n + b_n t + c_n t^2$  through the last rectangle at  $i = n$  we get

$$f(n) = \int_{n-\frac{1}{2}}^{n+\frac{1}{2}} (a_n + b_n t + c_n t^2) dt,$$

$$\frac{f(n) + f(n-1)}{2} = a_n + b_n (n - \frac{1}{2}) + c_n (n - \frac{1}{2})^2,$$

$$b_n + 2 c_n n + c_n = 0,$$

and for the constants

$$a_n = \frac{3}{4} (n^2 + n - \frac{1}{12}) f(n-1) - \frac{3}{4} (n^2 + n - \frac{1}{12}) f(n),$$

$$b_n = -\frac{3}{4} (1 + 2n) f(n-1) + \frac{3}{4} (1 + 2n) f(n),$$

$$c_n = \frac{3}{4} f(n-1) - \frac{3}{4} f(n).$$

Having obtained the constants for the graduating curve we will determine the moments of this curve in terms of those of the given frequency distribution.

*Notation:* Let the class interval be  $c = 1$ ; let  $\nu_s = \sum_{i=-m}^n i^s f(i)$  be the uncorrected  $s^{\text{th}}$  moment of the given frequency distribution about the given origin; let  $\mu_s = \sum_{i=-m}^n (i - \nu_1)^s f(i)$  be the uncorrected  $s^{\text{th}}$  moment of the given frequency distribution about its uncorrected mean; let  $\bar{\nu}_s$  be the corrected value of the  $s^{\text{th}}$  moment about the given origin; and let  $\bar{\mu}_s$  be the corrected value of the  $s^{\text{th}}$  moment about the corrected mean. Thus  $\nu_s$  and  $\mu_s$  apply to the rectangles, and  $\bar{\nu}_s$  and  $\bar{\mu}_s$  apply to the curves as follows:

$$\bar{\nu}_s = \sum_{i=-m+1}^{n-1} \int_{i-\frac{1}{2}}^{i+\frac{1}{2}} t^s (a + bt + ct^2) dt + \int_{-m-\frac{1}{2}}^{-m+\frac{1}{2}} t^s (a_1 + b_1 t + c_1 t^2) dt$$

$$+ \int_{n-\frac{1}{2}}^{n+\frac{1}{2}} t^s (a_n + b_n t + c_n t^2) dt,$$

$$\bar{\mu}_s = \sum_{i=-m+1}^{n-1} \int_{i-\frac{1}{2}}^{i+\frac{1}{2}} (t - \bar{\nu}_1)^s (a + bt + ct^2) dt + \int_{-m-\frac{1}{2}}^{-m+\frac{1}{2}} (t - \bar{\nu}_1)^s (a_1 + b_1 t + c_1 t^2) dt$$

$$+ \int_{n-\frac{1}{2}}^{n+\frac{1}{2}} (t - \bar{\nu}_1)^s (a_n + b_n t + c_n t^2) dt.$$

Using these symbols we have for the first moment about the given origin:

$$\bar{\nu}_1 = \sum_{i=-m+1}^{n-1} \int_{i-\frac{1}{2}}^{i+\frac{1}{2}} t (a + bt + ct^2) dt + \int_{-m-\frac{1}{2}}^{-m+\frac{1}{2}} t (a_1 + b_1 t + c_1 t^2) dt$$

$$+ \int_{n-\frac{1}{2}}^{n+\frac{1}{2}} t (a_n + b_n t + c_n t^2) dt$$



$$\begin{aligned}
&= \sum_{-m+1}^{n-1} \left[ ai + b \left( i^2 + \frac{1}{12} \right) + c \left( i^3 + \frac{i}{4} \right) \right] \\
&+ \left[ -a_1 m + b_1 \left( m^2 + \frac{1}{12} \right) - c_1 \left( m^3 + \frac{m}{4} \right) \right] \\
&+ \left[ a_n n + b_n \left( n^2 + \frac{1}{12} \right) + c_n \left( n^3 + \frac{n}{4} \right) \right].
\end{aligned}$$

Substituting the values for the constants this becomes

$$\begin{aligned}
\bar{v}_1 &= \sum_{-m+1}^{n-1} \left\{ i \left[ \left( \frac{5}{4} - 3i^2 \right) f(i) + \left( \frac{3i^2}{2} - \frac{i}{2} - \frac{1}{8} \right) f(i+1) \right. \right. \\
&\quad \left. \left. + \left( \frac{3i^2}{2} + \frac{i}{2} - \frac{1}{8} \right) f(i-1) \right] \right. \\
&+ (i^2 + \frac{1}{12}) [6if(i) + (\frac{1}{2} - 3i) f(i+1) - (\frac{1}{2} + 3i) f(i-1)] \\
&+ \left( i^3 + \frac{i}{4} \right) [-3f(i) + \frac{3}{2} f(i+1) + \frac{3}{2} f(i-1)] \left. \right\} \\
&+ \{ -m [\frac{3}{4} (m^2 + m - \frac{1}{12}) f(-m+1) - \frac{3}{4} (m^2 + m - \frac{1}{12}) f(-m)] \\
&+ (m^2 + \frac{1}{12}) [\frac{3}{4} (2m+1) f(-m+1) - \frac{3}{4} (2m+1) f(-m)] \\
&- \left( m^3 + \frac{m}{4} \right) [\frac{3}{4} f(-m+1) - \frac{3}{4} f(-m)] \} \\
&+ \{ n [\frac{3}{4} (n^2 + n - \frac{1}{12}) f(n-1) - \frac{3}{4} (n^2 + n - \frac{1}{12}) f(n)] \\
&+ \left( n^2 + \frac{1}{12} \right) \left[ -\frac{3}{4} (1+2n) f(n-1) + \frac{3}{4} (1+2n) f(n) \right] \\
&+ \left( n^3 + \frac{n}{4} \right) \left[ \frac{3}{4} f(n-1) - \frac{3}{4} f(n) \right] \}. \\
\bar{v}_1 &= \sum_{-m+1}^{n-1} \left[ if(i) + \frac{1}{24} f(i+1) - \frac{1}{24} f(i-1) \right] + \frac{1}{16} f(-m+1) \\
&\quad - \left( m + \frac{1}{16} \right) f(-m) - \frac{1}{16} f(n-1) + \left( n + \frac{1}{16} \right) f(n). \\
\sum_{-m+1}^{n-1} if(i) &= \sum_{-m}^n if(i) - (-m) f(-m) - n f(n) = v_1 + m f(-m) - n f(n). \\
\frac{1}{24} \sum_{i=-m+1}^{n-1} f(i+1) &= \frac{1}{24} \sum_{k=-m+2}^n f(k) = \frac{1}{24} \sum_{-m}^n f(k) - \frac{1}{24} f(-m+1) - \frac{1}{24} f(-m) \\
&= \frac{1}{24} - \frac{1}{24} f(-m+1) - \frac{1}{24} f(-m).
\end{aligned}$$



$$\begin{aligned}\frac{1}{24} \sum_{i=-m+1}^{n-1} f(i-1) &= \frac{1}{24} \sum_{j=-m}^{n-2} f(j) = \frac{1}{24} \sum_{-m}^n f(j) - \frac{1}{24} f(n-1) - \frac{1}{24} f(n) \\ &= \frac{1}{24} - \frac{1}{24} f(n-1) - \frac{1}{24} f(n).\end{aligned}$$

$$\begin{aligned}\bar{\nu}_1 &= \nu_1 + mf(-m) - nf(n) + \frac{1}{24} - \frac{1}{24} f(-m+1) - \frac{1}{24} f(-m) - \frac{1}{24} \\ &\quad + \frac{1}{24} f(n-1) + \frac{1}{24} f(n) + \frac{1}{16} f(-m+1) \\ &\quad - \left(m + \frac{1}{16}\right) f(-m) - \frac{1}{16} f(n-1) + \left(n + \frac{1}{16}\right) f(n).\end{aligned}$$

$$\bar{\nu}_1 = \nu_1 - \frac{5}{48} f(-m) + \frac{5}{48} f(n) + \frac{1}{48} f(-m+1) - \frac{1}{48} f(n-1).$$

Using this same notation and method for the higher moments we get

$$\begin{aligned}\bar{\mu}_2 &= \nu_2 - \frac{1}{12} - \bar{\nu}_1^2 + \left(\frac{5m}{24} + \frac{7}{80}\right) f(-m) + \left(\frac{5n}{24} + \frac{7}{80}\right) f(n) \\ &\quad + \left(\frac{-m}{24} - \frac{1}{240}\right) f(-m+1) + \left(\frac{-n}{24} - \frac{1}{240}\right) f(n-1).\end{aligned}$$

$$\begin{aligned}\bar{\mu}_3 &= \nu_3 - 3\bar{\nu}_1\bar{\mu}_2 - \frac{\bar{\nu}_1}{4} - \bar{\nu}_1^3 + f(-m) \left[\frac{-5}{16} m^2 - \frac{21}{80} m - \frac{17}{120}\right] \\ &\quad + f(n) \left[\frac{5}{16} n^2 + \frac{21}{80} n + \frac{17}{120}\right] + f(-m+1) \left[\frac{m^2}{16} + \frac{m}{80} + \frac{1}{120}\right] \\ &\quad + f(n-1) \left[\frac{-n^2}{16} - \frac{n}{80} - \frac{1}{120}\right].\end{aligned}$$

$$\begin{aligned}\bar{\mu}_4 &= \nu_4 - 4\bar{\mu}_3\bar{\nu}_1 - 6\bar{\mu}_2\bar{\nu}_1^2 - \bar{\nu}_1^4 - \frac{\bar{\mu}_2}{2} - \frac{\bar{\nu}_1^2}{2} - \frac{17}{64} \\ &\quad + f(-m) \left[\frac{5m^3}{12} + \frac{21m^2}{40} + \frac{17m}{30} + \frac{313}{1680}\right] + f(n) \left[\frac{5n^3}{12} + \frac{21n^2}{40} + \frac{17n}{30} + \frac{313}{1680}\right] \\ &\quad + f(-m+1) \left[\frac{-m^3}{12} - \frac{m^2}{40} - \frac{m}{30} - \frac{1}{336}\right] + f(n-1) \left[\frac{-n^3}{12} - \frac{n^2}{40} - \frac{n}{30} - \frac{1}{336}\right].\end{aligned}$$

#### SPECIAL CASES

The above formulae are rather long and in practice the special cases below will frequently be preferred.

(a) We may usually take the origin at or very near the middle of the range so that  $m = n$ , at least approximately.

If  $m = n$ :

$$\bar{\nu}_1 = \nu_1 - \frac{5}{48}f(-m) + \frac{5}{48}f(n) + \frac{1}{48}f(-m+1) - \frac{1}{48}f(n-1).$$

$$\begin{aligned}\bar{\mu}_2 = \nu_2 - \frac{1}{12} - \bar{\nu}_1^2 + \left(\frac{5m}{24} + \frac{7}{80}\right)[f(-m) + f(n)] \\ + \left(\frac{-m}{24} - \frac{1}{240}\right)[f(-m+1) + f(n-1)].\end{aligned}$$

$$\begin{aligned}\bar{\mu}_3 = \nu_3 - 3\bar{\nu}_1\bar{\mu}_2 - \frac{\bar{\nu}_1}{4} - \bar{\nu}_1^3 + \left[\frac{5m^2}{16} + \frac{21m}{80} + \frac{17}{120}\right][f(n) - f(-m)] \\ + \left[\frac{m^2}{16} + \frac{m}{80} + \frac{1}{120}\right][f(-m+1) - f(n-1)].\end{aligned}$$

$$\begin{aligned}\bar{\mu}_4 = \nu_4 - 4\bar{\mu}_3\bar{\nu}_1 - 6\bar{\mu}_2\bar{\nu}_1^2 - \bar{\nu}_1^4 - \frac{\bar{\mu}_2}{2} - \frac{\bar{\nu}_1^2}{2} - \frac{17}{64} \\ + \left[\frac{-m^3}{12} - \frac{m^2}{40} - \frac{m}{30} - \frac{1}{336}\right][f(-m+1) + f(n-1)] \\ + \left[\frac{5m^3}{12} + \frac{21m^2}{40} + \frac{17m}{30} + \frac{313}{1680}\right][f(-m) + f(n)].\end{aligned}$$

(b) Except in the abrupt cases the end frequencies and the difference between those next to the ends will be so small (relative to unity) that they will have a negligible effect on the corrections. If  $m = n$  as in (a), and if also

$$f(-m) = f(n) = 0 \text{ and } f(-m+1) - f(n-1) = 0:$$

$$\bar{\nu}_1 = \nu_1.$$

$$\bar{\mu}_2 = \nu_2 - \bar{\nu}_1^2 - \frac{1}{12} + f(-m+1)\left[\frac{-m}{12} - \frac{1}{120}\right].$$

$$\bar{\mu}_3 = \nu_3 - 3\bar{\nu}_1\bar{\mu}_2 - \frac{\bar{\nu}_1}{4} - \bar{\nu}_1^3.$$

$$\begin{aligned}\bar{\mu}_4 = \nu_4 - 4\bar{\mu}_3\bar{\nu}_1 - 6\bar{\mu}_2\bar{\nu}_1^2 - \bar{\nu}_1^4 - \frac{\bar{\mu}_2}{2} - \frac{\bar{\nu}_1^2}{2} - \frac{17}{64} \\ + f(-m+1)\left[\frac{-m^3}{6} - \frac{m^2}{20} - \frac{m}{15} - \frac{1}{168}\right].\end{aligned}$$

These formulae have been written in the form which makes the computing simple. The following makes a comparison with Sheppard's corrections easy.

$$\bar{\nu}_1 = \nu_1.$$

$$\bar{\mu}_2 = \mu_2 - \frac{1}{12} + f(-m+1) \left[ \frac{-m}{12} - \frac{1}{120} \right].$$

$$\bar{\mu}_3 = \mu_3 + \nu_1 \left( \frac{m}{4} + \frac{1}{40} \right) f(-m+1).$$

$$\bar{\mu}_4 = \mu_4 - \frac{\mu_2}{2} - \frac{43}{192} + f(-m+1) \left[ \frac{-m^3}{6} - \frac{m^2}{20} - \frac{m}{40} - \frac{1}{560} - \frac{m\nu_1^2}{2} - \frac{\nu_1^2}{20} \right].$$

The following special case is also useful in comparing my formulae with Sheppard's.

(c) Let  $f(-m) = \frac{1}{5} f(-m+1)$  and  $f(n) = \frac{1}{5} f(n-1)$ . This produces a graduating curve which is exactly tangent to the  $t$ -axis at the ends of the range and is everywhere continuous—though it does not have continuous derivatives at certain isolated points. It is, however, a curve which to the eye cannot be distinguished from the type assumed in the Euler-MacLaurin theorem, which lies at the base of Sheppard's formulae. My corrections become:

$$\bar{\nu}_1 = \nu_1,$$

$$\bar{\mu}_2 = \mu_2 - \frac{1}{12} + \frac{1}{15} [f(-m) + f(n)],$$

$$\bar{\mu}_3 = \mu_3 - \frac{\nu_1}{5} [f(-m) + f(n)] + \left[ \frac{-m}{5} - \frac{1}{10} \right] f(-m) + \left[ \frac{n}{5} + \frac{1}{10} \right] f(n),$$

$$\begin{aligned} \bar{\mu}_4 = \mu_4 - \frac{\mu_2}{2} - \frac{43}{192} + \frac{2}{5} \left[ (\nu_1 + m)^2 + \nu_1 + m + \frac{29}{84} \right] f(-m) \\ + \frac{2}{5} \left[ (\nu_1 - n)^2 - \nu_1 + n + \frac{29}{84} \right] f(n). \end{aligned}$$

Sheppard's are:

$$\bar{\nu}_1 = \nu_1,$$

$$\bar{\mu}_2 = \mu_2 - \frac{1}{12},$$

$$\bar{\mu}_3 = \mu_3,$$

$$\bar{\mu}_4 = \mu_4 - \frac{\mu_2}{2} + \frac{7}{240}.$$

Let us compare my results with Sheppard's in the very special case in which  $f(-m) = f(n) = 1/7$ ,  $f(0) = 5/7$ ,  $m = n = 1$ . The odd moments vanish. My corrections for  $\mu_2$  and  $\mu_4$  are

$$\bar{\mu}_2 = 0.2214, \quad \bar{\mu}_4 = 0.1870.$$

Sheppard's are

$$\bar{\mu}_2 = 0.2024, \quad \bar{\mu}_4 = 0.1720.$$

The numerical difference between the  $\bar{\mu}_2$ 's is 0.0190, and the numerical difference between the  $\bar{\mu}_4$ 's is 0.0150.

This example shows that Sheppard's corrections are not valid to the precision to which they are usually given if they are to be used for the purpose of correcting raw moments. The last term in the fourth moment correction,  $7/240$ , might equally well be, for example,  $-43/192$  as in my special case. This will become more evident to the reader if he will draw the curve indicated in this example. To the eye it will appear exactly like the kind specified in the Euler-MacLaurin theorem; for example, much like the normal curve. Now suppose one adopted for the moment the point of view (which I have criticized earlier) of starting with the curve used in this example, breaking it up into three partial areas and then finding the relation between the true and the raw moments. The partial areas found would be exactly those used in this example and this method would give us Sheppard's corrections, but they would not be exactly correct, for in this instance my formulae give exactly the relationship between the true and the raw moments. The difference is due to the fact that in this instance the assumptions permitting the use of the Euler-MacLaurin theorem in abbreviated form are not justified for this curve. But there is no way of telling at the outset, if one has given initially only the partial areas, whether precisely this curve or another which to the eye would appear very much like it is truly the curve which will graduate the same material when subjected to a finer classification.

## THE POINT BINOMIAL AND PROBABILITY PAPER

BY FRANK H. BYRON<sup>1</sup>

1. An approximation to the sum of a number of consecutive terms of the point binomial may be found graphically and quite expeditiously by means of so-called "probability paper." This paper is ruled so that the  $(x, y)$  graph of the equation of the integral of the normal curve

$$y = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{x^2}{2}} dx \quad (1)$$

is a straight line. Let the successive terms of the point binomial be represented as follows:

$$(p + q)^n = u_0 + u_1 + \cdots + u_t + \cdots + u_n, \quad (2)$$

where  $u_t = {}_nC_t p^{n-t} q^t$  and  $p \geq q$ . Then the  $(x, y)$  graph of the equation,

$$y = \sum_{i=0}^t u_i, \quad t + \frac{1}{2} = x, \quad (3)$$

*i.e.*, of the sum of first  $(t + 1)$  terms of this point binomial, is, in all but extreme cases, a set of points lying on a gently turning curve, so gently that its form may be represented closely by two straight lines, each passing through the median point as will be explained in the next section. As paper of this sort is readily obtainable, and as this method yields as great accuracy as is really useful in many problems, it is suggested that its use ought to be quite general.

2. **Sheppard's Corrections.** The formulae for the moments of the point binomial, mean =  $qn$ ,  $\sigma^2 = pqn$ , are exact without any corrections such as are used for grouped material. This fact has led us all (apparently) to assume that in fitting the curve to the point binomial one would get a better fit by equating the moments of the curve to the uncorrected moments of the point binomial rather than to the corrected moments. The studies made in connection with the preparation of this paper show that when the purpose is to equate areas to sums of terms the corrected moments should be used. The theoretical basis for this conclusion is as follows:

To simplify the argument let us suppose that one were seeking that curve of Charlier type,

$$F(x) = c_0\phi_0(x) + c_1\phi_1(x) + \cdots c_4\phi_4(x), \quad (4)$$

---

<sup>1</sup> With the assistance of Burton H. Camp.

(where  $\phi_0$  is the normal curve and  $\phi_1, \phi_2, \dots$  its successive derivatives) whose integral would best fit the graph of (3). Since fitting is required only at the isolated points  $x = \frac{1}{2}, 1\frac{1}{2}, 2\frac{1}{2}, \dots$ , it is clear that one might obtain this by the two following steps. First let  $f(x)$  be any function whose integral meets exactly the requirement at these isolated points. What values this integral has at other points does not for the moment concern us. There are an infinite number of such  $f(x)$  curves. Next let the  $c$ 's of (4) be so chosen that  $F(x)$  will fit  $f(x)$  as nearly as possible. The ordinary derivation of the  $c$ 's supposes that the fit between  $f(x)$  and  $F(x)$  is to be made by least squares, the residuals being weighted by the factor  $1/\sqrt{\phi(x)}$ . No matter what  $f(x)$  is chosen, the  $c$ 's can be determined so that the weighted integral of  $(f(x) - F(x))^2$  will be a minimum, but the value of this minimum will vary from one  $f(x)$  to another. We now desire to select that  $f(x)$  which will make this minimum value as small as possible, and it is reasonable to suppose that our best selection will be some  $f(x)$  which is as kindred to the nature of  $F(x)$  as possible. We shall not therefore choose an  $f(x)$  which oscillates wildly between the points where perfect fitting is required, (Fig. 1) nor yet an  $f(x)$  which is made up of the top bases of the point binomial

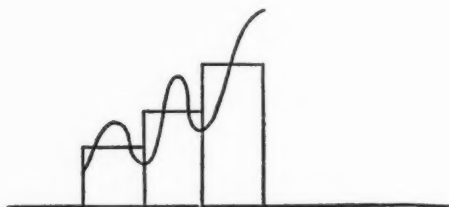


FIG. 1



FIG. 2

histogram; we shall prefer a modification (Fig. 2) of that histogram by a smoothing process. Such an  $f(x)$  will not have the exact moments of the point binomial, but, more nearly, those moments corrected for grouping. Then the determination of the  $c$ 's will come out in terms of these corrected moments, not in terms of the uncorrected moments. (In fact the uncorrected moments would be the exact moments of an  $f(x)$  having an oscillatory character between the important points.)

Of course, when  $n$  is large, the difference is too small to be noticed and the use of Sheppard's corrections is not worth while, and since  $n$  usually is large when approximations of this sort are needed, the point is not usually important. It was important in the computation of the tables of §4. Moreover, the use of Sheppard's corrections does not invariably yield better results, the gain being masked sometimes by other effects to be considered in §3. An excellent illustration of uniformly better results is in fitting  $(\frac{1}{2} + \frac{1}{2})^9$  by a curve of Type 4. The errors in the sums as derived from (4) with and without the corrections, is given on the following page.

$t$	0	1	2	3	4	5	6	7	8	9
With Corrections	.0002	.0001	-.0003	-.0001	.0000	.0001	.0003	-.0001	-.0002	.0000
Without Corrections	.0007	.0022	.0039	.0036	.0000	-.0036	-.0039	-.0022	-.0007	-.0001

3. **The Stubby End.** The other effects which mask this improvement are especially noticeable at the stubby end of a point binomial. We have to keep in mind here that the approximating curve (such as (4)), is required to turn a sharp corner, for, due to the least square method of fitting, it is just as important that it be close to zero when  $t$  is negative, as it is that it be close to  $u_0, u_1, \dots$  when  $t$  is positive. Therefore, in order to turn this corner it has to dip below the  $x$ -axis in the neighborhood of  $t = -\frac{1}{2}$ . This makes the approximating curve too low just to the right of  $t = -\frac{1}{2}$ , unless the whole curve be arbitrarily widened. This arbitrary widening is customarily performed by not using Sheppard's correction for  $\sigma$ , and the result is a betterment of the fit at these points but a corresponding loss over the rest of the infinite interval. A good example<sup>2</sup> is  $(\frac{2}{3} + \frac{1}{3})^{25}$ . The fit is worse at the left end when Sheppard's corrections are used but better over the rest of the interval.

The same difficulty arises in another connection. If we compare the closeness of fit to a point binomial made by  $F(x)$  as written in (4) and by  $F(x)$  as it would be written if  $c_4$  were zero, it often happens (as is well known) that the latter is actually slightly better on the average. How can this be true if the  $c$ 's are chosen by the method of least squares and the best choice as thus indicated makes  $c_4$  different from zero? The answer is that the  $c$ 's are chosen so that the fit is best over the infinite interval, not merely over the interval from  $t = -\frac{1}{2}$  to  $t = n + \frac{1}{2}$ , and that furthermore the distant points are weighted more heavily than those near the center. Thus it might happen that a choice, other than the least square choice, and one in which  $c_4$  would be zero, might be better for the restricted interval covered by the point binomial. This does happen especially when due to the abruptness of the stubby end of a very skew binomial, the curve has to dip below the axis in order to get by a sharp corner. A good example is the problem considered by Fry:<sup>3</sup>  $(\frac{9}{10} + \frac{1}{10})^{100}$ . All the effects mentioned are present here. The fit is on the average a little worse if  $c_4$  is not equal to zero over the point binomial interval, a little better over the infinite interval.

4. For graphical purposes a sufficiently good approximation to the median of  $(p + q)^n$ , is given by

$$M = nq - (p - q)/6.$$

<sup>2</sup> The true values are given on page 220 of Mathematical Part of Elementary Statistics, by Camp, D. C. Heath and Company, 1931.

<sup>3</sup> T. C. Fry, Probability and its Engineering Uses, p. 258, Van Nostrand, 1928.



The following tables enable us to find the first quartile  $Q_1$ , and the ninth decile  $D_9$ . The accuracy to which they can be plotted is only about one-tenth that to which they are given here. Therefore accurate interpolation is seldom necessary. The values of  $S_{t+1}$  are to be read from the graph at the points  $t + \frac{1}{2}$ , as indicated in the directions preceding the tables. The graphical method will be found efficient if one uses common sense in the computation. Numbers which are to be plotted should not be computed to a higher degree of accuracy than can be used graphically. In reading the values of  $S_{t+1}$  it is well to remember that the true values lie on a curve, and that outside the interval from  $Q_1$  to  $D_9$ , they are slightly less than those given by the straight line. Once the graph has

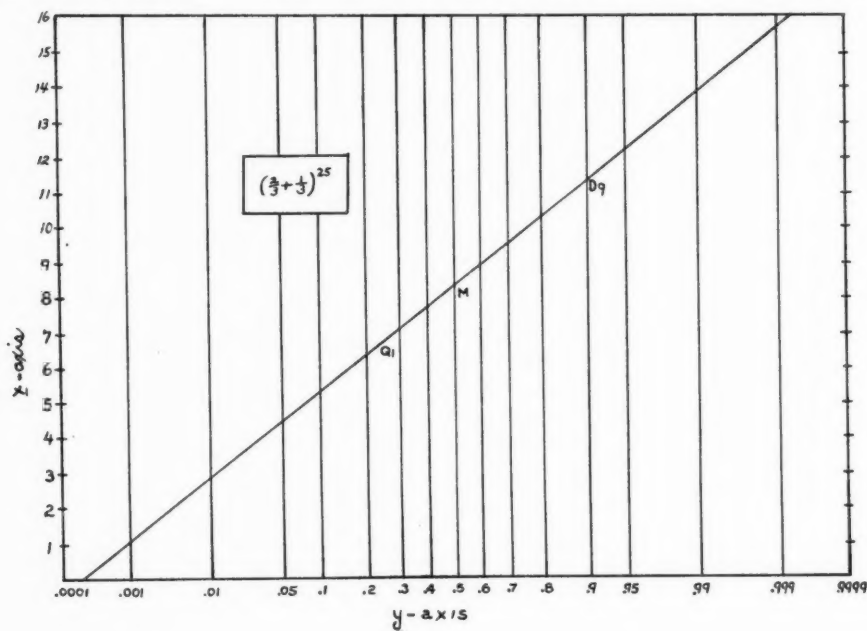


FIG. 3

been made, all the values of  $S_{t+1}$  can be read quickly; it is not necessary to make a separate computation for each  $t$ . This method is therefore specially advantageous when one wishes to find several sums of this sort for the same point binomial. It should also be noticed that one can tell from the appearance of the graph about how far the true sum would be from the two straight lines and so estimate the error to which his reading is liable.

**5. Illustration.** Find the sum of the first 7 terms of  $(\frac{2}{3} + \frac{1}{3})^{25}$ .

Here  $t = 6$ ,  $M = 8.278$ ,  $Q_1 = 6.726$ ,  $D_9 = 11.369$ . The graph shows that  $\sum_0^t = 0.224$ . The true value is 0.222. So the error is 0.002.





Values of  $x_2$ 

$n \backslash p$	2000	1000	750	500	400	300	200	100	75	50	25
.99	1.307	1.318	1.325	1.336	1.344	1.356	1.378	1.439	1.481	See Auxiliary Tables	
.98	299	307	311	318	323	330	343	376	396		
.97	295	301	304	310	314	319	329	353	367		
.96	293	298	301	306	309	313	321	341	352		
.95	292	296	299	303	305	309	316	332	342		
.94	291	295	297	300	303	306	312	327	335		
.93	290	293	295	298	301	304	309	322	329	1.342	1.374
.92	289	292	294	297	299	302	307	318	325	336	365
.91	289	292	293	296	298	300	305	315	321	331	357
.90	288	291	292	295	296	299	303	313	318	325	351
.88	287	290	291	293	295	297	300	309	313	321	341
.85	286	288	289	291	292	294	297	304	308	314	330
.80	285	287	288	289	290	291	293	298	301	306	317
.75	284	285	286	287	288	289	291	294	297	300	308
.70	284	285	285	286	286	287	288	291	293	295	301
.65	283	284	284	285	285	286	286	288	290	292	296
.60	283	283	283	284	284	284	285	286	287	288	291
.50	282	282	282	282	282	282	282	282	282	282	282

Auxiliary Table

$n \backslash p$	60	50	40	35	30	25	20
.99	1.525	1.575	1.663	1.740	1.871	2.149	3.209
.98	416	435	455	488	520	1.568	1.652
.97	381	394	413	433	445	472	514
.96	362	372	387	397	410	428	457
.95	350	359	370	378	389	405	425
.94	336	349	359	366	375	387	405

## INEQUALITIES AMONG AVERAGES

BY NILAN NORRIS

Numerous inequalities among averages of various types are condensed in the monotonic character of the function

$$\phi(t) = \left( \frac{x_1^t + x_2^t + \cdots + x_n^t}{n} \right)^{\frac{1}{t}}$$

of the positive numbers  $x_1, x_2, \dots, x_n$ , not all equal each to each. For  $t = -1$  this function is the harmonic mean; for  $t = 0$  it is the geometric mean; for  $t = 1$  the arithmetic mean; and for  $t = 2$  the root mean square. The relations among these four means which customarily are proved by special and disconnected methods appear easily as applications of the theorem that  $\phi(t)$  is an increasing function of  $t$ . That is, for any values of  $t_1$  and  $t_2$  such that  $-\infty < t_1 < t_2 < +\infty$ , it will be true that  $\phi(t_1) < \phi(t_2)$ . Several proofs of this theorem have been published, many of them very complex. An extremely simple proof is herewith presented.<sup>1</sup>

That  $\phi(t)$ ,  $\phi'(t)$  and  $\phi''(t)$  all exist and are continuous for all real values of  $t$  may be shown by expanding each of the quantities  $x_i^t$  in a series of powers of  $t$  and considering the remainders after each of the first three terms. The ordinary rule for evaluating forms reducing to  $0/0$ , which requires the function under consideration to be continuous and to have at least a continuous first derivative for  $t = 0$ , may then be applied to  $[\log \phi(t)]/t$  to show that  $\phi(0)$  is the geometric mean. It is clear that  $\phi(-\infty)$  and  $\phi(+\infty)$  are respectively the least and the greatest of the  $x_i$ . This fact and the monotonic property of  $\phi(t)$  make it evident that for each real value of  $t$ , the function may be regarded as an average in the usual sense that it lies within the range of the observations.

For a simple demonstration of the increasing character of  $\phi(t)$ , consider the auxiliary function

$$F(t) = t^2 \frac{\phi'(t)}{\phi(t)} = t^2 \frac{d}{dt} \left\{ \frac{1}{t} \log \frac{\sum x^t}{n} \right\} = t \frac{\sum x^t \log x}{\sum x^t} - \log \frac{\sum x^t}{n}.$$

It is clear that  $\phi'(t)$  has the same sign as  $F(t)$ . The theorem will be proved by showing that the sign of  $F(t)$  is positive for all values of  $t$  except zero, when  $\phi'(t)$  vanishes.

---

<sup>1</sup> Professor Harold Hotelling rendered invaluable assistance in condensing for publication the material herein presented from a more extended study of generalized mean value functions.

Differentiating the last expression with respect to  $t$ , one obtains upon simplification

$$F'(t) = \frac{t}{(\sum x^t)^2} [(\sum x^t) (\sum x^t \log^2 x) - (\sum x^t \log x)^2].$$

By Cauchy's inequality (known as Schwarz' inequality when applied to integrals instead of sums), the expression in square brackets is positive. Hence  $F'(t)$  has the same sign as  $t$ . Consequently  $F(t)$ , since it diminishes for negative values of  $t$  and increases for positive values, has a minimum for  $t = 0$ . But by direct substitution,  $F(0) = 0$ . It follows that  $F(t)$  and  $\phi'(t)$  are positive for all values of  $t$  other than zero. Therefore  $\phi(t)$  is an increasing function.

By direct general methods it is possible to show that

$$\phi'(0) = (\Pi x)^{\frac{1}{n}} \frac{1}{2n^2} [n \sum (\log x)^2 - (\sum \log x)^2].$$

This expression obviously vanishes only when  $n \sum (\log x)^2 = (\sum \log x)^2$ , a condition which is satisfied only in the trivial case when  $x_1 = x_2 = \dots = x_n$ .

A proof exactly parallel to that given above may be applied to integrals or, more generally, to Stieltjes integrals. The monotonic increasing character of  $\left[ \int_{x=0}^{\infty} x^t d\psi(x) \right]^{\frac{1}{t}}$  appears in this way if one assumes that  $\psi(x)$  is a non-decreasing function integrable in the Riemann-Stieltjes sense, such that  $\psi(\infty) - \psi(0) = 1$ , and such that  $\int_{x=0}^{\infty} x^t d\psi(x)$  exists for every real value of  $t$ . In terms of statistical theory, this consideration extends the theorem from samples to populations of a very general character.

Proof of the increasing character of  $\phi(t)$  has also been derived from Hölder's inequality, the demonstration being expressed in terms of Stieltjes integrals.<sup>2</sup> The simplest general proof of the monotonic attribute of  $\phi(t)$  heretofore published appears to be that of Paul Lévy.<sup>3</sup> As early as 1840 Bienaymé<sup>4</sup> presented a generalized form of  $\phi(t)$ , namely,

$$\left( \frac{c_1 a_1^m + c_2 a_2^m + \dots + c_n a_n^m}{c_1 + c_2 + \dots + c_n} \right)^{\frac{1}{m}},$$

and announced, without proof, its increasing character. In 1858 a proof of the monotonic quality of  $\phi(t)$  for special cases was published by Schlömilch.<sup>5</sup> Of

<sup>2</sup> J. Shohat, "Stieltjes Integrals in Mathematical Statistics," *Annals of Mathematical Statistics* (American Statistical Association, Ann Arbor, 1930), Vol. 1, No. 1, p. 84.

<sup>3</sup> *Calcul des Probabilités* (Gauthier-Villars et Cie., Paris, 1925), pp. 157 f.

<sup>4</sup> Jules Bienaymé, *Société Philomatique de Paris*, Extraits des Procès-Verbaux des Seances Pedant L'Anée 1840 (Imprimerie D'A. René et Cie., Paris, 1841), Seance du 13 juin 1840, p. 68.

<sup>5</sup> O. Schlömilch, "Ueber Mittelgrößen verschiedener Ordnungen," *Zeitschrift für Mathematik und Physik* (B. G. Teubner, Leipzig, 1858), Vol. 3, pp. 303 f.

the more recent general proofs of the increasing character of  $\phi(t)$  which have appeared, those of Jensen,<sup>6</sup> Pólya,<sup>7</sup> Jessen,<sup>8</sup> and Carathéodory<sup>9</sup> may be mentioned. A recent application of  $\phi(t)$  to index number theory is that of Professor John B. Canning.<sup>10</sup>

VASSAR COLLEGE.

---

<sup>6</sup> J. L. W. V. Jensen, "Sur Les Fonctions Convexes Et Les Inegalités Entre Les Valeurs Moyennes," *Acta Mathematica* (Beijers Bokförlagsaktielbolag, Stockholm, 1905), Vol. 30, pp. 183-185.

<sup>7</sup> G. Pólya and G. Szegő, *Aufgaben und Lehrsätze Aus Der Analysis* (Julius Springer, Berlin, 1925), Vol. I, pp. 54 f. and 210.

<sup>8</sup> Børge Jessen, "Bemaerkninger om koveskse Funktioner og Uligheder imellem Middellaerdier," *Matematisk Tidsskrift* (Charles Johansens Bogtrykkeri, Copenhagen, 1931), No. 2, 1931, pp. 26-28.

<sup>9</sup> Attributed to Professor Constantin Carathéodory in an unpublished manuscript of Professor Harold Hotelling.

<sup>10</sup> "A Theorem Concerning a Certain Family of Averages of a Certain Type of Frequency Distribution," a paper presented before a joint meeting of the American Statistical Association and the Econometric Society at Berkeley, California, June 22, 1934.

## MATHEMATICAL EXPECTATION OF PRODUCT MOMENTS OF SAMPLES DRAWN FROM A SET OF INFINITE POPULATIONS

BY HYMAN M. FELDMAN<sup>1</sup>

### Introduction

In the second part of his investigations, "On the Mathematical Expectation of Moments of Frequency Distributions,"<sup>2</sup> Tchouproff presented a method which may be interpreted as sampling from a set of infinite univariate populations. In the present paper this method is extended to the study of moments of product moments of samples drawn from a set of infinite bivariate populations. It is also shown how this method may be extended to populations of higher order by deriving some of the simpler formulae for populations of three and four variables.

Tchouproff's method has been criticised<sup>3</sup> because of the complicated algebra. On close examination it is found, however, that it is not the algebra which is complicated but rather the symbolism. Tchouproff introduced a great variety of symbols both in his derivations and in his results. As a consequence his work seems very intricate. If, however, the number of symbols is reduced, and the symbols themselves are simplified, which can be easily accomplished, the underlying idea of Tchouproff's method is found to be very simple.

Quite a complete study of product moments of any bivariate population has been made by Joseph Pepper in his "Studies in the Theory of Sampling."<sup>4</sup> His method is essentially an extension of Church's<sup>5</sup> method, in his studies of univariate populations, to bivariate populations. He does not, however, derive any generalized formulae. In the present study generalized formulae for both the first moment and the variance of product moments of any order are obtained.

It may be noted here, that all of Pepper's formulae for any infinite population can be obtained from those of the present study as special cases, by assuming that all the populations in the set are identical.

---

<sup>1</sup> A dissertation presented to the Board of Graduate Studies of Washington University in partial fulfillment of the requirements for the degree of Doctor of Philosophy, June 1933.

<sup>2</sup> *Biometrika*, Vol. XXI, Dec. 1929, pp. 231-258.

<sup>3</sup> Church, A. E. R. "On the Means and Squared Standard Deviations of Small Samples from any Population," *Biometrika*, Vol. XVIII, Nov., 1926, pp. 321-394.

<sup>4</sup> *Biometrika*, Vol. XXI, Dec. 1929, pp. 231-258.

<sup>5</sup> Church, A. E. R., "On the Means and Squared Standard Deviations of Small Samples from any Population," *Biometrika*, Vol. XVIII, Nov., 1926, pp. 321-394.

# CHAPTER I. Notations and Definitions

Let  $(X_1, Y_1), (X_2, Y_2), \dots (X_n, Y_n)$  be  $n$  bivariate populations each following any law of distribution whatever. The product moment of order  $a$  in  $X$  and  $b$  in  $Y$  of the  $k^{\text{th}}$  population will be denoted by  $P_{ab}^k$ . It is defined as

$$P_{ab}^k = E(X_k - a_k)^a (Y_k - b_k)^b \quad (1.11)$$

where

$$a_k = E(X_k), \quad b_k = E(Y_k), \quad (1.12)$$

and where the symbol  $E$  signifies the expected value or the mathematical expectation of a quantity.

Regarding each of the  $n$  populations of the set as infinite,<sup>6</sup> samples of  $n$  are drawn, each member of a sample from one of the  $n$  populations.<sup>7</sup> The individual which is drawn from the  $k^{\text{th}}$  population will be denoted by  $(x_k, y_k)$ ; and the product moment of order  $a$  in  $x$  and  $b$  in  $y$ , of such a sample will be denoted by  $p_{ab}$ . This product moment may then be defined as

$$p_{ab} = n^{-1} S (x_k - x)^a (y_k - y)^b \quad (1.13)$$

where

$$x = n^{-1} S x_k, \quad y = n^{-1} S y_k. \quad (1.14)$$

The symbols  $a$  and  $b$  will now be defined by the equations

$$a = n^{-1} S a_k, \quad b = n^{-1} S b_k. \quad (1.15)$$

$$\text{Obviously } E(x) = E(n^{-1} S x_k) = n^{-1} S E(X_k) = n^{-1} S a_k = a. \quad (1.16)$$

Similarly  $E(y) = b$ . That is, the mathematical expectation of the mean, of such a sample as was described above, is equal to the average of the means of all the populations.<sup>8</sup>

In order to make the equations as compact as possible the following additional symbols will be employed:

$$\begin{aligned} x_k - a_k &= u_k, & x - a &= u, & \text{and } u_k - u &= U_k \\ y_k - b_k &= v_k, & y - b &= v, & \text{and } v_k - v &= V_k \end{aligned} \quad (1.17)$$

also  $a_k - a = A_k, b_k - b = B_k$ .

From the above definitions it easily follows that

$$E(u_k) = E(v_k) = E(U_k) = E(V_k) = E(u) = E(v) = 0. \quad (1.18)$$

<sup>6</sup> The term infinite is used here in the probability sense. It is defined very clearly by Church in his "Means and Squared Standard Deviations of Small Samples," *Biometrika*, Vol. XVIII, Nov., 1926, p. 322.

<sup>7</sup> It may be easily shown that this is equivalent to drawing a sample of  $n$  from a set of any finite number of populations. The number drawn from each population, however, must be specified. See *Biometrika*, Vol. XIII, 1920-21, p. 295, footnote.

<sup>8</sup> This, of course, is a result of the Lexis Theory, for Poisson and Lexis Series.



The notation is now completed with the definition of the symbol  $Q_{ij}$  by the equation:

$$Q_{ij} = S(a_k - a)^i (b_k - b)^j = SA_k^i B_k^j. \quad (1.19)$$

## CHAPTER II. The Mathematical Expectation of $p_{ab}$

The mathematical expectation of  $p_{ab}$  will be denoted by  $\bar{p}_{ab}$ . In the terminology of moments this would be called the mean or first moment of the distribution of  $p_{ab}$ .

**1. The Mathematical Expectation of  $p_{11}$ .** According to the above notation the expected value of  $p_{11}$  is  $\bar{p}_{11}$ . By definition

$$\bar{p}_{11} = E(p_{11}) = En^{-1}S(x_i - x)(y_i - y), \quad (2.11)$$

and obviously  $En^{-1}S(x_i - x)(y_i - y) = n^{-1}SE(x_i - x)(y_i - y)$ .

Writing

$$x_i - x = [(x_i - a_i) - (x - a)] + [a_i - a] = U_i + A_i$$

$$y_i - y = [(y_i - b_i) - (y - b)] + [b_i - b] = V_i + B_i,$$

equation (2.11) may be written as

$$\begin{aligned} \bar{p}_{11} &= n^{-1}SE(U_i + A_i)(V_i + B_i) \\ &= n^{-1}SE(U_i V_i) + n^{-1}SA_i E(V_i) + n^{-1}SB_i E(U_i) + n^{-1}SE(A_i B_i). \end{aligned}$$

Since for any given population  $A_i$  and  $B_i$  are constants, it follows that  $E(A_i B_i) = A_i B_i$ . Hence

$$n^{-1}SE(A_i B_i) = n^{-1}SA_i B_i = n^{-1}Q_{11}.$$

Making use of (1.18), it is seen that the terms  $n^{-1}SA_i E(V_i)$  and  $n^{-1}SB_i E(U_i)$  are zero. The only term left to evaluate is therefore  $n^{-1}SE(U_i V_i)$ . Since  $U_i$  and  $V_i$  are symmetric functions of the corresponding small letters, their product is symmetric in  $u, v$ . There is therefore no loss in generality if attention is concentrated on a single subscript, say 1.

We may therefore write

$$n^{-1}SE(U_i V_i) = n^{-1}E(U_1 V_1) + n^{-1}SE(U_i V_i)_{*2}$$

Remembering that  $U_i = u_i - u = u_i - n^{-1}Su_i$ , we may write,

$$\begin{aligned} U_i &= u_i - u = u_i - n^{-1}(u_1 + u_2 + \cdots + u_n) \\ &= n^{-1}[n_1 u_i - (u_1 + u_2 + \cdots + u_{i-1} + u_{i+1} + \cdots + u_n)] \end{aligned}$$

---

\* The 2 at the bottom of the  $S$  simply indicates that the summation begins with  $i = 2$ .



where  $n_1 = n - 1$ . In general,  $n_i$  will denote the number  $n - i$ . Similarly

$$V_i = n^{-1}[n_1 v_1 - (v_1 + v_2 + \dots + v_{i-1} + v_{i+1} + \dots + v_n)].$$

Thus

$$\begin{aligned} n^{-1}SE(U_i V_i) &= n^{-3}E(n_1 u_1 - u_2 - \dots - u_n)(n_1 v_1 - v_2 - \dots - v_n) \\ &+ n^{-3}SE(n_1 u_i - u_1 - \dots - u_{i-1} - u_{i+1} - \dots - u_n) \\ &\quad (n_1 v_i - v_1 - \dots - v_{i-1} - v_{i+1} - \dots - v_n). \end{aligned}$$

When the right hand side of the last equation is expanded the only terms which appear are of the form  $E(u_i v_i)$  and  $E(u_i v_j)$ . The last one must vanish for  $u_i$  and  $v_j$  are independent and hence  $E(u_i v_j) = E(u_i)E(v_j) = 0$ . From the last equation above it is easily seen that the coefficient of  $E(u_i v_i)$  is

$$n^{-3}(n_1^2 + n_1) = n^{-3} n_1(n_1 + 1) = n^{-2} n_1;$$

and because of the symmetry this is obviously the coefficient of any term of that form. Hence

$$n^{-1}SE(U_i V_i) = n^{-2} n_1 SE(u_i v_i).$$

Since  $u_i = x_i - a_i$ ,  $v_i = y_i - b_i$ , then

$$E(u_i v_i) = E(x_i - a_i)(y_i - b_i) = E(X_i - a_i)(Y_i - b_i) = P_{11}^i$$

and in general,

$$E(u_k^i v_k^j) = P_{ij}^k. \quad (2.12)$$

We thus get the formula

$$\bar{p}_{11} = n^{-2} n_1 SP_{11}^i + n^{-1} Q_{11}. \quad (1)$$

Now suppose all the  $n$  populations are identical. Then all the  $A$ 's and also all the  $B$ 's vanish and therefore,  $Q_{11} = 0$ . The formula (1) thus becomes

$$\bar{p}_{11} = \frac{n-1}{n} P_{11}. \quad (1')$$

This is exactly Pepper's formula for  $\bar{p}_{11}$  for an infinite population.<sup>9</sup>

## 2. The Mathematical Expectation of $p_{21}$ . By definition

$$\bar{p}_{21} = En^{-1}S(x_i - \bar{x})^2(y_i - \bar{y}). \quad (2.21)$$

<sup>9</sup> *Biometrika*, Vol. XXI, p. 233, Eq. A,  $N = \infty$ . As was already stated in the introduction, all the formulae of the present study reduce to Pepper's when the above assumption is made.

Proceeding as above it is seen that

$$\begin{aligned} En^{-1}S(x_i - x)^2(y_i - y) &= n^{-1}SE(x_i - x)^2(y_i - y) \\ &= n^{-1}SE(U_i + A_i)^2(V_i + B_i) = n^{-1}SE(U_i^2 V_i) + 2n^{-1}SE(U_i V_i A_i) \\ &\quad + n^{-1}SE(U_i^2 B_i) + n^{-1}SE(V_i A_i^2) + 2n^{-1}SE(A_i B_i U_i) + n^{-1}SE(A_i^2 B_i) \dots (2.22) \end{aligned}$$

It is quite evident that the two terms before the last vanish. To evaluate the remaining terms, we employ the reasoning of section 1 of this chapter and write:

$$\begin{aligned} SE(U_i^2 V_i) &= E(U_i^2 V_i) + SE(U_i^2 V_i) \\ &= n^{-3}E(n_1^* u_1 - u_2 - \dots)(n_1 v_1 - v_2 - \dots) + n^{-3}SE(n_1 u_i - u_1 - \dots) \\ &\quad (n_1 v_i - v_1 - \dots). \end{aligned}$$

Since terms of the form  $E(u_i^2 v_i)$  vanish, only the coefficient of the term  $E(u_i^2 v_i)$  must be found. Again considering the subscript 1, the coefficient of  $E(u_1^2 v_1)$  is easily found from the last equation to be

$$n^{-3}(n_1^3 - n_1) = n^{-3}n_1(n_1 + 1)(n_1 - 1) = n^{-2}n_1 n_2.$$

Thus

$$n^{-1}SE(U_i^2 V_i) = n^{-2}n_1 n_2 SE(u_i^2 v_i) = n^{-2}n_1 n_2 SP_{21}^i. \quad (2.23)$$

For the second term of (2.22) we have

$$\begin{aligned} SE(U_i V_i A_i) &= E(U_i V_i A_i) + SE(U_i V_i A_i) \\ &= n^{-2}E(n_1 u_1 - u_2 - \dots)(n_1 v_1 - v_2 - \dots)A_i + n^{-2}SE(n_1 u_i - u_1 - \dots) \\ &\quad (u_1 v_i - v_1 - \dots)A_i. \end{aligned}$$

The coefficient of  $E(u_1 v_1)$  in the first term of the right hand side of the last equation is  $n^{-2}n_1^2 A_1$ . In the second term it is  $n^{-2}SA_i = -n^{-2}A_1$ , since  $SA_i = 0$ .

It therefore follows that

$$2n^{-1}SE(U_i V_i A_i) = 2n^{-2}n_2 SP_{11}^i A_i. \quad (2.24)$$

Quite similarly

$$n^{-1}SE(U_i^2 B_i) = n^{-2}n_2 SP_{20}^i B_i, \quad (2.25)$$

and it is obvious that

$$n^{-1}SE(A_i^2 B_i) = n^{-1}Q_{21}. \quad (2.26)$$

---

\* Note that the  $u$  which has the coefficient  $n_1$  does not occur among the  $u$ 's which have the negative sign.

We thus get the formula

$$\bar{p}_{21} = n^{-3}n_1n_2SP_{21}^i + n^{-2}n_2S(2P_{11}^iA_i + P_{20}^iB_i) + n^{-1}Q_{21}. \quad (2)$$

### 3. The Mathematical Expectation of $p_{31}$ and $p_{22}$ .

$$\begin{aligned} \bar{p}_{31} &= En^{-1}S(x_i - x)^3(y_i - y) = n^{-1}SE(x_i - x)^3(y_i - y) \\ &= n^{-1}SE(U_i + A_i)^3(V_i + B_i) = n^{-1}S\{E(U_i^3V_i + U_i^3B_i + 3U_i^2V_iA_i \\ &\quad + 3U_i^2A_iB_i + 3U_iV_iA_i^2 + 3U_iA_i^2B_i + V_iA_i^3 + A_i^3B_i)\}. \end{aligned} \quad (2.31)$$

The two terms before the last are zero. The last term is

$$n^{-1}SE(A_i^3B_i) = n^{-1}Q_{31}. \quad (2.32)$$

By (2.23) and (2.24) and some slight manipulation

$$\begin{aligned} &3n^{-1}SE(U_i^2A_iB_i + U_iV_iA_i^2) \\ &= 3n^{-3}n_2S(P_{20}^iA_iB_i + P_{11}^iA_i^2) + 3n^{-3}(Q_{11}SP_{20}^i + Q_{20}SP_{11}^i), \end{aligned} \quad (2.33)$$

and by (2.22)

$$n^{-1}SE(U_i^3B_i + 3U_i^2V_iA_i) = n^{-4}(n_1^3 + 1)S(P_{30}^iB_i + 3P_{21}^iA_i). \quad (2.34)$$

The only new term which is to be evaluated is  $SE(U_i^3V_i)$ . This may be written as follows:

$$SE(U_i^3V_i) = n^{-4}SE(n_1u_i - u_1 - \dots)^3(n_1v_i - v_1 - \dots).$$

When the right hand side is expanded it is found that the only non-vanishing terms are of the form  $E(U_i^3V_i)$  and  $E(u_i^2u_jv_i)$ . Only two subscripts, therefore, have to be considered. Without any loss in generality these may be taken as 1 and 2, and the right hand side of the last equation may then be written as follows:

$$\begin{aligned} SE(n_1u_i - u_1 - \dots)^3(n_1v_i - v_1 - \dots) &= E(n_1u_1 - u_2 - \dots)^3(n_1v_1 - v_2 - \dots) \\ &\quad + E(n_1u_2 - u_1 - \dots)^3(n_1v_1 - v_2 - \dots) + SE(n_1u_i - u_1 - u_2 - \dots)^3 \\ &\quad \quad \quad (n_1v_i - v_1 - v_2 - \dots). \end{aligned}$$

From this last expansion it is easily seen that the coefficient of  $E(u_i^3v_i)$  is  $(n_1^4 + n_1)$  and that of  $E(u_i^2u_jv_i)$ ,  $(6n_1^2 + 3n_2) = 3(2n_1^2 + n_2)$ . We thus finally obtain

$$SE(U_i^3V_i) = n^{-4}\{(n_1^4 + n_1)SE(u_i^3v_i) + 3(2n_1^2 + n_2)SE(u_i^2u_jv_i)\}.$$

But by (2.12)  $E(u_i^3v_i) = P_{31}^i$ , and since  $u_i$  and  $u_j$  and  $u_i$  and  $v_j$  are independent  $E(u_i^2u_jv_i) = E(u_i^2)E(u_jv_i) = P_{20}^iP_{11}^i$ . Whence

$$E(U_i^3V_i) = n^{-4}\{(n_1^4 + n_1)SP_{31}^i + 3(2n_1^2 + n_2)SP_{20}^iP_{11}^i\}. \quad (2.35)$$

From (2.31) and the succeeding equations we finally get

$$\begin{aligned}\bar{p}_{31} = & n^{-5}\{(n_1^4 + n_1)SP_{31}^i + 3(2n_1^2 + n_2)SP_{20}^{*i}P_{11}^i\} \\ & + n^{-4}\{(n_1^3 + 1)S(P_{30}^iB_i + 3P_{21}^iA_i)\} + 3n^{-3}\{(n_1^2 - 1)S(P_{20}^iA_iB_i + P_{11}^iA_i^2) \\ & + Q_{11}SP_{20}^i + Q_{20}SP_{11}^i\} + n^{-1}Q_{31}.\end{aligned}\quad (3)$$

The derivation of  $\bar{p}_{22}$  is so similar to that of  $\bar{p}_{31}$ , that it would be mere repetition to go through the details again. We shall therefore merely write down the formula for  $\bar{p}_{22}$  which is

$$\begin{aligned}\bar{p}_{22} = & n^{-5}\{(n_1^4 + n_1)SP_{22}^i + (2n_1^2 + n_2)S(P_{20}^iP_{02}^i + 4P_{11}^iP_{11}^i)\} \\ & + 2n^{-4}\{(n_1^3 + 1)S(P_{21}^iB_i + P_{12}^iA_i)\} + n^{-3}\{(n_1^2 - 1)S(P_{20}^iB_i^2 + 4P_{11}^iA_iB_i \\ & + P_{02}^iA_i^2) + Q_{20}SP_{02}^i + Q_{02}SP_{20}^i + 4Q_{11}SP_{11}^i\} + n^{-1}Q_{22}.\end{aligned}\quad (4)$$

#### 4. The Mathematical Expectation of the General Product Moment $p_{ab}$ .

So far, formulae for the mathematical expectation of  $p_{ab}$ , for particular values of  $a$  and  $b$ , have been derived. The method used in deriving these is, however, perfectly general, and now, that it has been sufficiently illustrated, it can be easily generalized.

By definition we have

$$\bar{p}_{ab} = E[n^{-1}S(x_i - x)^a(y_i - y)^b].$$

Making use of the notation of Chapter I this may be written as

$$n\bar{p}_{ab} = ES(U_i + A_i)^a(V_i + B_i)^b = \sum_{q,r=0}^{a,b} C_q^a C_r^b SE(U_i^{a-q} V_i^{b-r} A_i^q B_i^r) \quad (2.41)$$

where

$$C_q^a = \frac{a!}{q!(a-q)!}, \quad C_r^b = \frac{b!}{r!(b-r)!}.$$

Expressing the  $U$ 's and  $V$ 's in terms of the  $u$ 's and  $v$ 's and setting  $a - q = l$ ,  $b - r = m$ ; we may write for a particular pair of values  $q$  and  $r$ :

$$n^{l+m} SE(U_i^l V_i^m A_i^q B_i^r) = SE(n_1 u_i - u_1 - \dots)^l (n_1 v_i - v_1 - \dots)^m A_i^q B_i^r. \quad (2.42)$$

Consider, now, the general term in the expansion of the right hand side of (2.42). It is of the form:

$$\frac{l!m!}{\Pi\alpha_k! \Pi\beta_k!} (-1)^{l+m} (-n_1)^{\alpha_k + \beta_k} E(n_1 u_{i1}^{\alpha_1} \dots u_{ik}^{\alpha_k} v_{i1}^{\beta_1} \dots v_{ik}^{\beta_k} A_i^q B_i^r), \quad (2.43)$$

where  $\Pi\alpha_k! = \alpha_1! \alpha_2! \dots \alpha_k!$

\* In this case, and also in the formulae that follow, whenever two or more indices appear in a summation, it will be understood that no two of them can have the same value simultaneously.

For particular sets of values  $j_1, j_2, \dots, j_k, \alpha_1, \alpha_2, \dots, \alpha_k$ , and  $\beta_1, \beta_2, \dots, \beta_k$ , this term will appear in every member of the summation of the right hand side of (2.42), and its coefficient will differ only in the exponent of  $(-n_1)$  and in the subscript  $i$  of  $A^q B^r$ . Because of the symmetry there is no loss in generality if we take for  $j_1, j_2, \dots, j_k$ , the first  $k$  integers. We now break up the summation of the right hand side of (2.42) as follows:

$$\begin{aligned} & \sum_1^n SE(n_1 u_i - u_1 - \dots)^l (n_1 v_i - v_1 - \dots)^m A_i^q B_i^r \\ &= E(n_1 u_1 - u_2 - \dots)^l (n_1 v_1 - v_2 - \dots)^m A_1^q B_1^r \\ &+ E(n_1 u_2 - u_1 - \dots)^l (n_1 v_2 - v_1 - \dots)^m A_2^q B_2^r + \dots + E(n_1 u_k - u_1 - \dots)^l \\ & (n_1 v_k - v_1 - \dots)^m A_k^q B_k^r + \sum_{i=k+1}^n E(n_1 u_i - u_1 - \dots)^l \\ & (n_1 v_i - v_1 - \dots)^m A_i^q B_i^r. \end{aligned} \quad (2.44)$$

From (2.44) we easily get for the total coefficient (excluding the numerical factor) the expression

$$\sum_{h=1}^k (-n_1)^{\alpha_h + \beta_h} A_h^q B_h^r + \sum_{h=k+1}^n A_h^q B_h^r.$$

Writing

$$\sum_{k+1}^n A_h^q B_h^r = \sum_1^n A_h^q B_h^r - \sum_1^k A_h^q B_h^r = Q_{qr} - \sum_1^k A_h^q B_h^r,$$

the general term, (2.43), together with the total coefficient, may then be written as

$$(-1)^{l+m} \frac{l! m!}{\Pi \alpha_h! \Pi \beta_h!} \left\{ \sum_{h=1}^k [(-n_1)^{\alpha_h + \beta_h} - 1] A_h^q B_h^r + Q_{qr} \right\} \sum_{h=1}^k E u_h^{\alpha_h} v_h^{\beta_h}.$$

Since  $u_i$  and  $u_j, v_i$  and  $v_j$ , and  $u_i$  and  $v_j$  are independent while  $u_i$  and  $v_i$  are not, we have:

I.  $E \Pi u_h v_h = \Pi E u_h v_h = \Pi P_{\alpha_h \beta_h}^h$

II. Any term in which  $\alpha_h + \beta_h = 1$  must vanish.

From II it follows that the maximum number of subscripts which can appear in any term in the expansion of (2.42), i.e. the upper limit of  $k$ , which will be denoted by  $t$ , cannot exceed  $(l+m)/2$ . In fact when  $l+m$  is even,  $t = (l+m)/2$ , while when  $l+m$  is odd,  $t$  is the largest integer less than  $(l+m)/2$ .

Making use of (2.41), the equations following it, and the reasoning of the last paragraph, we finally get the formula:

$$\begin{aligned} n(-n)^{a+b} \tilde{p}_{ab} &= (a!) (b!) \sum_{j,h=1}^n \sum_{q,r=0}^{a,b} \frac{(-n)^{q+r}}{q! r!} \sum_{\alpha_h=0, \beta_h=0}^{a-q, b-r} S \quad S \\ & \left\{ \sum_{h=1}^k [(-n_1)^{\alpha_h + \beta_h} - 1] A_{j,h}^q B_{j,h}^r + Q_{qr} \right\} \Pi \frac{P_{\alpha_h \beta_h}^{j,h}}{\alpha_h! \beta_h!}. \end{aligned} \quad (5)$$

The following restrictions on the  $\alpha$ 's and  $\beta$ 's must be observed

$$(a) \alpha_1 + \alpha_2 + \dots + \alpha_k = a - q$$

$$(b) \beta_1 + \beta_2 + \dots + \beta_k = b - r$$

$$(c) \alpha_h + \beta_h \neq 1.$$

In case the  $n$  populations are identical (5) reduces as follows: For  $q = 0$ ,  $r = 0$ ,  $A_i^0 = 1$ ,  $B_i^0 = 1$ , and  $Q_{00} = n$ ; while in every other case  $A_i^q B_i^r = 0$ ,  $Q_{qr} = 0$ . The summations with respect to  $q$  and  $r$ , therefore disappear.

Consider now the summations

$$\sum_{j_1=1}^n \sum_{j_2=1}^n \dots \sum_{j_k=1}^n P_{j_1}^{i_1} P_{j_2}^{i_2} \dots P_{j_k}^{i_k}.$$

Since all the populations are the same we may drop the  $j$  by actually carrying out the indicated summations. If, then, there are  $c$  repetitions among the  $k$  pairs of integers  $\alpha_h \beta_h$ , in which  $\alpha_1 \beta_1$ ,  $\alpha_2 \beta_2$ ,  $\dots$   $\alpha_c \beta_c$ , are repeated  $l_1$ ,  $l_2$ ,  $\dots$   $l_c$  times respectively, then we have;

$$\sum_{j_1=1}^n \dots \sum_{j_k=1}^n P_{\alpha_h \beta_h}^{j_h} = \frac{k! C_k^n}{l_1! l_2! \dots l_c!} \Pi P_{\alpha_h \beta_h}.$$

We thus arrive at the following corollary: The mathematical expectation,  $\bar{p}_{ab}$ , of the product moment,  $p_{ab}$ , in samples of  $n$  from a single infinite population having any law of distribution is given by

$$n(-n)^{a+b} \bar{p}_{ab} = \sum_{\alpha_h, \beta_h=0}^{a,b} \frac{(a!)(b!)}{\Pi \alpha_h! \Pi \beta_h!} \sum_{k=1}^t \left[ \sum_{h=1}^n S(-n_1)^{\alpha_h + \beta_h} + n_k \right] \frac{k! C_k^n}{l_1! \dots l_c!} \Pi P_{\alpha_h \beta_h}^* \quad (5')$$

*Note:* In deriving these general formulae it was assumed that  $n > t$ . There is however, no loss in generality in this assumption. For, if  $t > n$ , we may suppose that,  $x_{n+1} = x_{n+2} = \dots = x_c = 0$ , and hence  $P_{\alpha\beta}^{n+1} = \dots = P_{\mu\nu}^t = 0$ , and thus the above reasoning is still valid.

**5. Formulae for  $\bar{p}_{41}$ ,  $\bar{p}_{32}$ ,  $\bar{p}_{51}$ ,  $\bar{p}_{42}$ ,  $\bar{p}_{33}$ .** Formulae for  $\bar{p}_{ab}$  in which  $a + b = 5, 6, 7, 8$  have been obtained. But for  $(a + b) > 6$  these formulae become very long, and since these will be of no use in the subsequent work, only those of order 5 and 6 are given below.

$$\begin{aligned} \bar{p}_{41} = & n^{-6} \{ (n_1^5 - n_1) S P_{41}^i + 2nn_2^2 S (2P_{30}^i P_{11}^i + 3P_{21}^i P_{20}^i) \} \\ & + n^{-5} \{ (n_1^4 + n_1) S (P_{40}^i B_i + 4P_{31}^i A_i) + 6nn_2 S (P_{20}^i B_i P_{20}^i + 2P_{11}^i A_i P_{20}^i) \} \end{aligned}$$

\* This is a generalization of Pepper's results for  $N = \infty$ . See *Biometrika* Vol. XXI, pp. 231-240.

† The symbol  $P_{11}^i A_i P_{20}^i$  is an abbreviation of the full term  $(A_i + A_j) (P_{11}^i P_{20}^i + P_{11}^j P_{20}^j)$ . Similar abbreviations will be used in the other formulae.

$$+ 2n^{-4} \{ (n_1^3 + 1) S(2P_{30}^i A_i B_i + 3P_{21}^i A_i^2) - 2Q_{11} S P_{30}^i - 3Q_{20} S P_{21}^i \} \\ + 2n^{-3} \{ (n_1^2 - 1) S(2P_{11}^i A_i^3 3P_{20}^i A_i^2 B_i) + 2Q_{30} S P_{11}^i + 3Q_{21} S P_{20}^i \} + n^{-1} Q_{41}. \quad (6)$$

$$\bar{p}_{32} = n^{-6} \{ (n_1^5 - n_1) S P_{32}^i + n n_2^2 S (P_{30}^i P_{02}^j + 6P_{21}^i P_{11}^j + 3P_{20}^i P_{12}^j) \} \\ + n^{-5} \{ (n_1^4 - 1) S(2P_{31}^i B_i + 3P_{22}^i A_i) + 3n n_2 S (P_{20}^i P_{11}^j B_i + [P_{20}^i P_{02}^j \\ + 4P_{11}^i P_{11}^j] A_i) \} + n^{-4} \{ (n_1^3 + 1) S(P_{30}^i B_i^2 + 6P_{21}^i A_i B_i \\ + 3P_{12}^i A_i^2) - Q_{02} S P_{30}^i - 6Q_{11} S P_{21}^i - 3Q_{20} S P_{12}^i \} + n^{-3} \{ (n_1^2 - 1) S(3P_{20}^i A_i B_i^2 \\ + 6P_{11} A_i B_i + P_{02}^i A_i^3) + 3Q_{12} S P_{20}^i + 6Q_{21} S P_{11}^i + Q_{30} S P_{02}^i \} + n^{-1} Q_{32}. \quad (7)$$

$$\bar{p}_{51} = n^{-7} \{ (n_1^6 + n_1) S P_{51}^i + 5(n_1^4 + n_1^2 + n_2) S (P_{40}^i P_{11}^j \\ + 2P_{31}^i P_{20}^j) - 10(2n_1^3 - n_2) S P_{21}^i P_{30}^j + 30(3n_1^2 + n_3) S P_{20}^i P_{20}^j P_{11}^k \} + n^{-6} \{ (n_1^5 \\ + 1) S (P_{50}^i B_i + 5P_{41}^i A_i) + 10(n_1^3 + 1) S^* [2P_{30}^i P_{20}^j B_i + (2P_{30}^i P_{11}^j \\ + 3P_{21}^i P_{20}^j) A_i] - 10n n_2 S^* [2P_{30}^i P_{20}^j B_i + (2P_{30}^i P_{11}^j + 3P_{21}^i P_{20}^j) A_i] \} \\ + 5n^{-5} \{ (n_1^4 - 1) S (P_{40}^i A_i B_i + 2P_{31}^i A_i^2) + 6n n_2 S (P_{20}^i P_{20}^j A_i B_i + 2P_{20}^i P_{11}^j A_i^2) \\ + Q_{11} S (P_{40}^i + 6P_{20}^i P_{20}^j) + 2Q_{20} S (P_{31}^i + 6P_{20}^i P_{11}^j) \} + 10n^{-4} \{ (n_1^3 \\ + 1) S (P_{30}^i A_i^2 B_i + P_{21}^i A_i^3) - Q_{21} S P_{30}^i - Q_{30} S P_{21}^i \} + 5n^{-3} \{ (n_1^3 - 1) S (2P_{20}^i A_i^3 B_i \\ + P_{11}^i A_i^4) + 2Q_{31} S P_{20}^i + Q_{40} S P_{11}^i \} + n^{-1} Q_{51}. \quad (8)$$

$$\bar{p}_{42} = n^{-7} \{ (n_1^6 + n_1) S P_{42}^i + (n_1^4 + n_1^2 + n_2) S (P_{40}^i P_{02}^j + 8P_{31}^i P_{11}^j \\ + 6P_{22}^i P_{20}^j) + 4(2n_1^3 - n_2) S (P_{30}^i P_{12}^j + 3P_{21}^i P_{21}^j) + 6(3n_1^2 + n_3) S (P_{20}^i P_{20}^j P_{02}^k \\ + 4P_{20}^i P_{11}^j P_{11}^k) \} + 2n^{-6} \{ (n_1^5 + 1) S (P_{41}^i B_i + 2P_{32}^i A_i) + 2(n_1^3 + 1) S [(2P_{30}^i P_{11}^j \\ + 3P_{21}^i P_{20}^j) B_i + (P_{30}^i P_{02}^j + 6P_{21}^i P_{11}^j + 3P_{12}^i P_{20}^j) A_i] - 2n n_2 S [(2P_{30}^i P_{11}^j \\ + 3P_{21}^i P_{20}^j) B_i + (P_{30}^i P_{12}^j + 6P_{21}^i P_{11}^j + 3P_{12}^i P_{20}^j) A_i] \} + n^{-5} \{ (n_1^4 - 1) S (P_{40}^i B_i^2 \\ + 8P_{31}^i A_i B_i + 6P_{22}^i A_i^2) + 6n n_2 S [P_{20}^i P_{21}^j B_i^2 + 4P_{20}^i P_{11}^j A_i B_i + (P_{20}^i P_{02}^j \\ + 4P_{11}^i P_{11}^j) A_i^2] + Q_{02} S (P_{40}^i + 6P_{20}^i P_{20}^j) + 8Q_{11} S (P_{31}^i + 3P_{20}^i P_{11}^j) \\ + 6P_{20} S (P_{22}^i + P_{20}^i P_{02}^j + 4P_{11}^i P_{11}^j) \} + 4n^{-4} \{ (n_1^3 + 1) S (P_{30}^i A_i B_i^2 \\ + 3P_{21}^i A_i^2 B_i + P_{12}^i A_i^3) - Q_{12} S P_{30}^i - 3Q_{21} S P_{21}^i + Q_{30} S P_{12}^i \} \\ + n^{-3} \{ S (6P_{20}^i A_i^2 B_i^2 + 8P_{11}^i A_i^3 B_i + P_{02}^i A_i^4) + Q_{40} S P_{02}^i + 8Q_{31} S P_{11}^i \\ + 6Q_{22} S P_{20}^i \} + n^{-1} Q_{42}. \quad (9)$$

$$\bar{p}_{33} = n^{-7} \{ (n_1^6 + n_1) S P_{33}^i + 3(n_1^4 + n_1^2 + n_2) S (P_{31}^i P_{22}^j + 3P_{22}^i P_{11}^j \\ + P_{13}^i P_{20}^j) - (2n_1^3 - n_2) S (P_{30}^i P_{03}^j + 9P_{21}^i P_{12}^j) + 9(3n_1^2 + n_3) S (P_{20}^i P_{11}^j P_{02}^k$$

\* The repetition of this expression signifies that  $A$  and  $B$  factors are coupled only with those  $P$  factors which have corresponding indices.



$$\begin{aligned}
& + 4P_{11}^i P_{11}^j P_{11}^k \} + 3n^{-6} \{ (n_1^5 + 1) S(P_{32}^i B_i + P_{23}^i A_i) + (n_1^3 + 1) S[(P_{30}^i P_{02}^j \\
& + 6P_{21}^i P_{11}^j + 3P_{12}^i P_{20}^j) B_i + (P_{03}^i P_{20}^j + 6P_{12}^i P_{11}^j + 3P_{21}^i P_{02}^j) A_i] \\
& - nn_2 S[(P_{30}^i P_{02}^j + 6P_{21}^i P_{11}^j + 3P_{12}^i P_{20}^j) B_i + (P_{03}^i P_{20}^j + 6P_{12}^i P_{11}^j \\
& + 3P_{21}^i P_{02}^j) A_i] \} + 3n^{-5} \{ (n_1^4 - 1) S(P_{31}^i B_i^2 + 3P_{22}^i A_i B_i + P_{13}^i A_i^2) \\
& + 3n_1 n_2 S[P_{20}^i P_{11}^j B_i^2 + (P_{20}^i P_{02}^j + 4P_{11}^i P_{11}^j) A_i B_i + P_{11}^i P_{02}^j A_i^2] \\
& + S[Q_{02}(P_{31}^i + 3P_{20}^i P_{11}^j) + 3Q_{11}(P_{22}^i + P_{20}^i P_{02}^j + 4P_{11}^i P_{11}^j) + Q_{20}(P_{13}^i \\
& + 3P_{02}^i P_{11}^j)] \} + n^{-4} \{ (n_1^3 + 1) S(P_{30}^i B_i^3 + 9P_{21}^i A_i B_i^2 + 9P_{12}^i A_i^2 A_i + P_{03}^i A_i^3) \\
& - S(Q_{03} P_{30}^i + 9Q_{12} P_{21}^i + 9Q_{21} P_{12}^i + Q_{30} P_{03}^i) \} + 3n^{-3} \{ (n_1^2 - 1) S(P_{20}^i A_i B_i^2 \\
& + 3P_{11}^i A_i^2 B_i^2 + P_{02}^i A_i^3 B_i) + S(Q_{13} P_{20}^i + 3Q_{22} P_{11}^i + Q_{31} P_{02}^i) \} + n^{-1} Q_{33}. \quad (10)
\end{aligned}$$

### CHAPTER III. The Mathematical Expectation of the Variance of $p_{ab}$

1. The Symbols  ${}_2m_{p_{ab}}$  and  ${}_2M_{p_{ab}}$ . Denoting the variance of  $p_{ab}$  by  $m$  and the mathematical expectation of  ${}_2m_{p_{ab}}$  by  ${}_2M_{p_{ab}}$ , we have the definition,

$$\begin{aligned}
{}_2m_{p_{ab}} &= \{ n^{-1} S(x_i - x)^a (y_i - y)^b - \bar{p}_{ab} \}^2 \\
&= n^{-2} S^2(x_i - x)^a (y_i - y)^b - 2n^{-1} \bar{p}_{ab} S(x_i - x)^a (y_i - y)^b + \bar{p}_{ab}^2, \text{ and} \\
{}_2M_{p_{ab}} &= E({}_2m_{p_{ab}}) = E\{ n^{-2} S^2(x_i - x)^a (y_i - y)^b - 2n^{-1} \bar{p}_{ab} S(x_i - x)^a (y_i - y)^b + \bar{p}_{ab}^2 \} \\
&= n^{-2} E[S(x_i - x)^{2a} (y_i - y)^{2b}] + 2n^{-2} E[S(x_i - x)^a (x_j - x)^a (y_i - y)^b (y_j - y)^b] \\
&\quad - 2n^{-1} \bar{p}_{ab} E[S(x_i - x)^a (y_i - y)^b] + \bar{p}_{ab}^2 = n^{-1} \bar{p}_{2a2b} \\
&\quad + 2n^{-2} E[S(x_i - x)^a (y_i - y)^b (x_j - x)^a (y_j - y)^b] - \bar{p}_{ab}^2. \quad (3.11)
\end{aligned}$$

Before attempting to expand the right hand side of (3.11) for any values  $a, b$  we shall derive the formula for  ${}_2M_{p_{11}}$  to illustrate the procedure.

2. The Mathematical Expectation of  ${}_2m_{p_{11}}$ . By (3.11) we have

$${}_2M_{p_{11}} = n^{-1} \bar{p}_{22} + 2n^{-2} E[S(x_i - x)(y_i - y)(x_j - x)(y_j - y)] - \bar{p}_{11}^2. \quad (3.21)$$

The first term is given by (4) and the last by (1). The only new term is the middle one. To expand it let us write it in terms of  $U$  and  $V$ . We then have:

$$\begin{aligned}
n^{-2} SE[(x_i - x)(y_i - y)(x_j - x)(y_j - y)] &= n^{-2} SE[(U_i + A_i)(V_i \\
&+ B_i)(U_j + A_j)(V_j + B_j)] = n^{-2} \{ SE[U_i V_i U_j V_j + (U_i V_i U_j B_j + U_j V_j U_i B_i) \\
&+ (U_i V_i V_j A_j + U_j V_j V_i A_i) + (U_i V_i A_j B_j + U_j V_j A_i B_i) \\
&+ (U_i V_j A_i B_i + U_j V_i A_j B_j) + U_i U_j B_i B_j + V_i V_j A_i A_j + 4 \text{ vanishing terms} \\
&+ A_i B_i A_j B_j \} \}. \quad (3.22)
\end{aligned}$$

The evaluation of the last term is very simple. For

$$SE(A_i B_i A_j B_j) = S(A_i B_i A_j B_j),$$

and from the elementary theory of symmetric functions we have:

$$S(A_i B_i A_j B_j) = \frac{S^2(A_i B_i) - S(A_i^2 B_i^2)}{2}.$$

Hence

$$SE(A_i B_i A_j B_j) = \frac{S^2(A_i B_i) - S(A_i^2 B_i^2)}{2} = \frac{Q_{11}^2 - Q_{22}}{2}. \quad (3.23)$$

To expand the first term and also the remaining ones, we return to the  $u, v$ , notation defined in Chapter I. We then write

$$SE(U_i V_i U_j V_j) = n^{-4} SE[(n_1 u_i - u_1 - \dots)(n_1 v_i - v_1 - \dots) \\ (n_1 u_j - u_1 - \dots)(n_1 v_j - v_1 - \dots)].$$

The only terms which can appear in the expansion of the right hand side of the last equation have the following form:

$$E(u_i^2 v_i^2), \quad E(u_i^2 v_j^2), \quad E(u_i v_i u_j v_j),$$

i.e., exactly those which appear in the evaluation of  $\bar{p}_{22}$ . Remembering the symmetry, there will be no loss in generality if we take for  $i$  and  $j$  the integers 1 and 2. To find the coefficients of the three characteristic terms, the above summation may be broken up as follows:

$$n^4 SE(U_i V_i U_j V_j) = E[(n_1 u_1 - u_2 - \dots)(n_1 v_1 - v_2 - \dots)(n_1 u_2 - u_1 - \dots) \\ (n_1 v_2 - v_1 - \dots)] + E\{[n_1 u_1 - u_2 - \dots](n_1 v_1 - v_2 - \dots) + (n_1 u_2 - u_1 - \dots) \\ (n_1 v_2 - v_1 - \dots)] S(n_1 u_i - u_1 - \dots)(n_1 v_i - v_1 - \dots)\} + SE[(n_1 u_i - u_1 - \dots) \\ (n_1 v_i - v_1 - \dots)(n_1 u_j - u_1 - \dots)(n_1 v_j - v_1 - \dots)]. \quad (3.24)$$

Writing the three terms in a row and their coefficients from the three parts of (3.24) in columns below these terms, we get the following scheme:

	$E(u_1^2 v_1^2)$	$E(u_1^2 v_2^2 + u_2^2 v_1^2)$	$E(u_1 v_1 u_2 v_2)$
	$n_1^2$	$n_1^2$	$(n_1^2 + 1)^2$
	$n_2(n_1^2 + 1)$	$-2n_1 n_2$	$2n_2^3$
	$\frac{n_2 n_3}{2}$	$\frac{n_2 n_3}{2}$	$2n_2 n_3$
Total	$\frac{nn_1(2n_1 - 1)}{2}$	$\frac{-nn_3}{2}$	$n(n_1^3 + n_1^2 - 3n_1 + 3)$
coeff.			

With the aid of the above equations we finally get:

$$SE(U_i V_i U_j V_j) = n^{-4} \left\{ \frac{n_1 n (2n_1 - 1)}{2} SP_{22}^i - \frac{nn_3}{2} SP_{20}^i P_{02}^j + n(nn_1^2 - 3n_2) SP_{11}^i P_{11}^j \right\}$$

Proceeding in the same way we find:

$$SE(U_i V_i U_j B_j + U_j V_j U_i B_i) = n^{-3} (2n_1^2 + n_2) SP_{21}^i B_i$$

$$SE(U_i V_i V_j A_j + U_j V_j V_i A_i) = n^{-3} (2n_1^2 + n_2) SP_{12}^i A_i$$

$$SE(U_i V_i A_j B_j + U_j V_j A_i B_i) = -nn_2 SP_{11}^i A_i B_i + (n_1^2 + n_2) Q_{11} SP_{11}^i$$

$$SE(U_i V_j A_j B_i + U_j V_i A_i B_j) = 2n SP_{11}^i - Q_{11} SP_{11}^i$$

$$SE(U_i U_j B_i B_j + V_i V_j A_i A_j) = nS(P_{20}^i B_i^2 + P_{02}^i A_i^2) - \frac{1}{2}S(Q_{20} P_{02}^i + Q_{02} P_{20}^i).$$

Collecting terms and simplifying we finally get:

$$\begin{aligned} {}_2M_{P_{11}} &= n^{-4} \{ n_1^2 SP_{22}^i + S(P_{20}^i P_{02}^j + 2P_{11}^i P_{11}^j) - n^2 S(P_{11}^i)^2 \} \\ &+ 2n^{-3} n_1 \{ S(P_{21}^i B_i + P_{12}^i A_i) \} + n^{-2} \{ S(P_{20}^i B_i^2 + 2P_{11}^i A_i B_i + P_{02}^i A_i^2) \}. \quad (11) \end{aligned}$$

Corollary 1. In case  $X_i = Y_i$ , i.e., when the set of populations are univariate, (11) becomes

$${}_2M_{P_{20}} = n^{-4} \{ n_1^2 S[P_{40}^i - (P_{20}^i)^2] + 4SP_{20}^i P_{20}^j \} + 4n^{-3} n_1 SP_{30}^i A_i + 4n^{-2} SP_{20}^i A_i^2. \quad (11')$$

This is Tchouproff's formula for the expected value of the variance of samples of  $n$ .<sup>10</sup>

Corollary 2. In case the  $n$  populations are identical (11) becomes

$${}_2M_{P_{11}} = n^{-3} n_1 [n_1 P_{22} + P_{20} P_{02} - n_2 P_{11}^2]. \quad (11'')$$

**3. The Mathematical Expectation of  ${}_2M_{P_{ab}}$ .** We now return to the general equation

$${}_2M_{P_{ab}} = n^{-1} \bar{p}_{2a2b} - \bar{p}_{ab}^2 + 2n^{-2} \sum_{i=1, j=1}^n E(x_i - x)^a (y_i - y)^b (x_j - x)^a (y_j - y)^b. \quad (3.11)$$

\* Since  $E(u_i^2 v_i^2) = P_{22}^i$ ,  $E(u_i^2 v_j^2) = P_{20}^i P_{02}^j$ , etc.

<sup>10</sup> See *Biometrika*, Vol. XIII p. 295.

<sup>11</sup> See *Biometrika*, Vol. XXI p. 234, Cor. 1.

The first two terms are given by (5). To evaluate the last term we write:

$$\begin{aligned}
 SE[(x_i - x)^a (y_i - y)^b (x_j - x)^a (y_j - y)^b] &= SE[(U_i + A_i)^a (V_i + B_i)^b (U_j + A_j)^a \\
 (V_j + B_j)^b] &= SE(U_i^a V_i^b U_j^a V_j^b) + \sum_{r_1, r_2, r_3, r_4=0}^{a, a, b, b,} S C_{r_1}^a C_{r_2}^a C_{r_3}^b C_{r_4}^b \\
 SE(U_i^a V_i^b U_j^a V_j^b A_i^{r_1} B_i^{r_2} A_j^{r_3} B_j^{r_4}) &= n^{-2(a+b)} SE\{(n_1 u_i - \dots)^a (n_1 v_i - \dots)^b \\
 (n_1 u_j - \dots)^a (n_1 v_j - \dots)^b + S_{r_1 \dots r_4} n^{(r_1+r_2+r_3+r_4)} C_{r_1}^a \dots C_{r_4}^b SE[(n_1 u_i - \dots)^a \\
 (n_1 v_i - \dots)^b (n_1 u_j - \dots)^a (n_1 v_j - \dots)^b] &+ A_i^{r_1} B_i^{r_2} A_j^{r_3} B_j^{r_4}, \quad (3.31)
 \end{aligned}$$

where  $\alpha = a - r_1$ ,  $\beta = a - r_2$ ,  $\gamma = b - r_3$ ,  $\delta = b - r_4$ .

The right hand side of (3.31) has been broken up into two parts because the first part is symmetrical, while the second part, in general, is not except when  $r_1 = r_2$ , and  $r_3 = r_4$ .

Let us now consider the expression

$$SE[(n_1 u_i - \dots)^a (n_1 v_i - \dots)^b (n_1 u_j - \dots)^a (n_1 v_j - \dots)^b]. \quad (3.32)$$

This is a double summation in which  $c_{ij} = c_{ji}$  and in which the diagonal terms,  $c_{ii}$ , are missing.

Consider next a general term of  $k$  factors from the expansion of each bracket of (3.32). As we are dealing with symmetric functions, there will be no loss in generality if we consider the first  $k$  subscripts only; and if we let the lower limits of the exponents of the  $u$ 's and  $v$ 's begin with zero we may consider that each parenthesis of a given bracket contributes exactly  $k$  factors. Such a term, omitting the coefficient, may be written as follows:

$$\begin{aligned}
 E(u_1^{\alpha_1} \dots u_k^{\alpha_k} v_1^{\beta_1} \dots v_k^{\beta_k} u_1^{\alpha'_1} \dots u_k^{\alpha'_k} v_1^{\beta'_1} \dots v_k^{\beta'_k}) &= \prod_{h=1}^k E(u_h^{\alpha_h + \alpha'_h} v_h^{\beta_h + \beta'_h}) \\
 &= \prod_{h=1}^k P^h(\alpha_h + \alpha'_h)(\beta_h + \beta'_h). \quad (3.33)
 \end{aligned}$$

This term occurs in every one of the  $\frac{1}{2}nn_1$  brackets of (3.32), having the same numerical coefficient in every one of them, which is

$$\frac{(a!)^2 (b!)^2}{\Pi \alpha_h! \Pi \alpha'_h! \Pi \beta_h! \Pi \beta'_h!}. \quad (3.34)$$

To obtain the  $n_1$  coefficient of (3.33) we break up (3.32) into the following partial summations:

$$\begin{aligned}
 E[(n_1 u_i - \dots)^a (n_1 v_i - \dots)^b (n_1 u_j - \dots)^a (n_1 v_j - \dots)^b] &= E[(n_1 u_1 - \dots)^a \\
 (n_1 v_1 - \dots)^b (n_1 u_2 - \dots)^a (n_1 v_2 - \dots)^b] &+ \dots + E[(n_1 u_{k-1} - \dots)^a
 \end{aligned}$$

$$\begin{aligned}
& (n_1 v_{k-1} - \dots)^b (n_1 u_k - \dots)^a (n_1 v_k - \dots)^b + \sum_{i=1}^k E \left[ (n_1 u_i - \dots)^a \right. \\
& \left. (n_1 v_i - \dots)^b \sum_{j=k+1}^n (n_1 u_j - \dots)^a (n_1 v_j - \dots)^b \right] + \sum_{i,j=k+1}^n [E \{ (n_1 u_i - \dots)^a \\
& (n_1 v_i - \dots)^b (n_1 u_j - \dots)^a (n_1 v_j - \dots)^b \}].
\end{aligned}$$

From this equation we get for the total coefficient in  $n$  of the term (3.33) the following expression:

$$\sum_{h,h'=1}^k (-n_1)^{\alpha_h + \alpha_{h'} + \beta_h + \beta_{h'}} + n_k \sum_{h=1}^k [(-n_1)^{\alpha_h + \beta_h} + (-n_1)^{\alpha_{h'} + \beta_{h'}}] + C_2^{n_k}.$$

The following restrictions on the  $\alpha$ 's and  $\beta$ 's must be observed.

$$\begin{aligned}
(a) \quad & \alpha_1 + \alpha_2 + \dots + \alpha_k = a \quad (b) \quad \beta_1 + \beta_2 + \dots + \beta_h = b \\
& \alpha'_1 + \alpha'_2 + \dots + \alpha'_k = a \quad \beta'_1 + \beta'_2 + \dots + \beta'_k = b \\
(c) \quad & \alpha_h + \alpha_{h'} + \beta_h + \beta_{h'} \neq 1.
\end{aligned}$$

From (c) we obtain the upper limit of  $k$ , namely:  $t = a + b$ .

Combining the various above equations we finally obtain:

$$\begin{aligned}
& (n)^{2(a+b)} S(U_i^a V_i^b U_j^a V_j^b) = (a!)^2 (b!)^2 \sum_{i,h=1}^n \sum_{\alpha_h, \alpha_{h'}, \beta_h, \beta_{h'}=0}^{n,b} S \\
& \sum_{k=1}^t \left\{ \sum_{h,h'=1}^k (-n_1)^{\alpha_h + \beta_h + \alpha_{h'} + \beta_{h'}} + n_k \sum_{h=1}^k [(-n_1)^{\alpha_h + \beta_h} + (-n_1)^{\alpha_{h'} + \beta_{h'}}] + C_2^{n_k} \right\} \\
& \frac{\Pi P_{(\alpha_h + \alpha_{h'}) (\beta_h + \beta_{h'})}^{J_h}}{\Pi \alpha_h! \Pi \alpha_{h'}! \Pi \beta_h! \Pi \beta_{h'}!}. \quad (3.35)
\end{aligned}$$

Turning to the second part of (3.31) let us consider the expression

$$\sum_{i=1, j=1}^n E[(n_1 u_i - \dots) (n_1 v_i - \dots) (n_1 u_j - \dots) (n_1 v_j - \dots) A_i^{r_1} B_i^{r_2} A_j^{r_3} B_j^{r_4}]$$

for a given set of  $r$ 's. The term (3.33) may also be considered as a general term of this last expression; of course, the exponents of the  $u$ 's and  $v$ 's will be different in this case. In order to evaluate the complete coefficient of a term like (3.33) we again write;

$$\begin{aligned}
& SE[(n_1 u_i - \dots)^a (n_1 v_i - \dots)^\gamma (n_1 u_j - \dots)^\delta (n_1 v_j - \dots)^\delta A_i^{r_1} B_i^{r_2} A_j^{r_3} B_j^{r_4} \\
& = E[(n_1 u_1 - \dots)^a (n_1 v_1 - \dots)^\gamma (n_1 u_2 - \dots)^\delta (n_1 v_2 - \dots)^\delta A_1^{r_1} B_1^{r_2} A_2^{r_3} B_2^{r_4}] \\
& + E[(n_1 u_2 - \dots)^a (n_1 v_2 - \dots)^\gamma (n_1 u_1 - \dots)^\delta (n_1 v_1 - \dots)^\delta A_2^{r_1} B_2^{r_2} A_1^{r_3} B_1^{r_4} \\
& + \dots + E[(n_1 u_k - \dots)^a (n_1 v_k - \dots)^\gamma (n_1 u_{k-1} - \dots)^\delta (n_1 v_{k-1} - \dots)^\delta]
\end{aligned}$$

$$\begin{aligned}
 & A_k^{r_1} B_k^{r_3} A_{k-1}^{r_2} B_{k-1}^{r_4} + \sum_{i=1}^k S E[(n_1 u_i - \dots)^\alpha (n_1 v_i - \dots)^\gamma A_i^{r_1} B_i^{r_3}] \sum_{j=k+1}^n \\
 & (n_1 u_j - \dots)^\beta (n_1 v_j - \dots)^\delta A_j^{r_2} B_j^{r_4} + \sum_{j=1}^k S E[(n_1 u_j - \dots)^\beta (n_1 v_j - \dots)^\delta \\
 & A_j^{r_2} B_j^{r_4}] \sum_{i=k+1}^n (n_1 u_i - \dots)^\alpha (n_1 v_i - \dots)^\gamma A_i^{r_1} B_i^{r_3} + \sum_{i,j=k+1}^n S E[(n_1 u_i - \dots)^\alpha \\
 & n_1 v_i - \dots)^\gamma (n_1 u_j - \dots)^\beta (n_1 v_j - \dots)^\delta A_i^{r_1} B_i^{r_3} A_j^{r_2} B_j^{r_4}]. \quad (3.36)
 \end{aligned}$$

It is now quite easy to write down the complete coefficient of a term of the form (3.33). The numerical coefficient of this term is the same in every bracket of (3.36), and is

$$\frac{(-1) S_i (a - r_1)! (a - r_2)! (b - r_3)! (b - r_4)!}{\Pi \alpha_h! \Pi \alpha'_h! \Pi \beta_h! \Pi \beta'_h!} \quad (3.37)$$

The coefficient in  $n_1$  and  $A_i^{r_1} B_i^{r_3} A_j^{r_2} B_j^{r_4}$  is broken up by (3.36) into the following four parts:

$$\text{I. } \sum_{h=1, h'=1}^k (-n_1)^{\alpha_h + \alpha'_h + \beta_h + \beta'_h} A_h^{r_1} B_h^{r_3} A_{h'}^{r_2} B_{h'}^{r_4}, \text{ from the first } k(k-1) \text{ brackets.}$$

$$\begin{aligned}
 \text{II. } \sum_{h=1}^k (-n_1)^{\alpha_h + \beta_h} A_h^{r_1} B_h^{r_3} \sum_{h'=k+1}^n A_{h'}^{r_2} B_{h'}^{r_4} &= \sum_{h=1}^k (-n_1)^{\alpha_h + \beta_h} A_h^{r_1} B_h^{r_3} \\
 &\left[ Q_{r_2 r_4} - \sum_{h'=1}^k A_{h'}^{r_2} B_{h'}^{r_4} \right],
 \end{aligned}$$

from the next  $k(n-k)$  brackets. Similarly

$$\text{III. } \sum_{h'=1}^k (-n_1)^{\alpha_{h'} + \beta_{h'}} A_{h'}^{r_2} B_{h'}^{r_4} \left[ Q_{r_1 r_3} - \sum_{h=1}^k A_h^{r_1} B_h^{r_3} \right], \text{ from the next } k(n-k).$$

And finally:

$$\begin{aligned}
 \text{IV. } \sum_{i,j=k+1}^n A_i^{r_1} B_i^{r_3} A_j^{r_2} B_j^{r_4} &= \sum_1^n A_h^{r_1} B_h^{r_3} \sum_1^n A_{h'}^{r_2} B_{h'}^{r_4} - \sum_1^n A_h^{(r_1+r_2)} B_h^{(r_3+r_4)} \\
 &- \sum_{h,h'=1}^k A_h^{r_1} B_h^{r_3} A_{h'}^{r_2} B_{h'}^{r_4} - \sum_{h=1}^n A_h^{r_1} B_h^{r_3} \sum_{h'=1}^k A_{h'}^{r_2} B_{h'}^{r_4} - \sum_{h=1}^n A_{h'}^{r_2} B_{h'}^{r_4} \sum_{h=1}^k A_h^{r_1} B_h^{r_3} \\
 &+ 2 \sum_{h=1}^k A_h^{r_1} B_h^{r_3} \sum_{h'=1}^k A_{h'}^{r_2} B_{h'}^{r_4} = Q_{r_1 r_3} Q_{r_2 r_4} - Q_{(r_1+r_2)(r_3+r_4)} - Q_{r_1 r_3} \sum_1^n A_h^{r_2} B_h^{r_4} \\
 &- Q_{r_2 r_4} \sum_1^n A_h^{r_1} B_h^{r_3} - \sum_{h,h'=1}^k A_h^{r_1} B_h^{r_3} A_{h'}^{r_2} B_{h'}^{r_4} + 2 \sum_{h=1}^k A_h^{r_1} B_h^{r_3} \sum_{h'=1}^k A_{h'}^{r_2} B_{h'}^{r_4}, \text{ from the}
 \end{aligned}$$

last  $c_2^{nk}$  brackets.

The restrictions on the  $\alpha$ 's and  $\beta$ 's differ from those given above in that  $a$  is replaced by  $a - r_1$  and  $a - r_2$ , and  $b$  by  $b - r_3$  and  $b - r_4$ ; and from the restriction (c) we get for the upper limit of  $k$ , in this case,

$$t_1 = \frac{\alpha + \beta + \gamma + \delta}{2} = a + b - \frac{r_1 + r_2 + r_3 + r_4}{2}$$

when  $Sr_i$  is even, or the greatest integer less than  $\frac{S\alpha}{2}$  when  $Sr_i$  is odd.

Combining (3.37) with  $C_{r_1}^a \dots C_{r_4}^b$  we get for the general numerical coefficient in the expansion of the second part of (3.31), the expression

$$\frac{(-1)^{Sr_i} (a!)^2 (b!)^2}{\Pi r_i! \Pi \alpha_h! \Pi \alpha'_h! \Pi \beta_h! \Pi \beta'_h!}.$$

By an obvious manipulation we have

$$\begin{aligned} \text{I} + \text{II} + \text{III} + \text{IV} &= \sum_{h,h'=1}^k \left[ (-n_1)^{\alpha_h + \beta_h + \alpha'_h + \beta'_h} - 1 \right] A_h^{r_1} B_h^{r_3} A_h^{r_2} B_h^{r_4} + Q_{r_2 r_4} \\ &\quad \sum_{h=1}^k \left[ (-n_1)^{\alpha_h + \beta_h} - 1 \right] A_h^{r_1} B_h^{r_3} + Q_{r_1 r_3} \sum_{h=1}^k \left[ (-n_1)^{\alpha'_h + \beta'_h} - 1 \right] A_h^{r_2} B_h^{r_4} \\ &\quad - \sum_{h=1}^k A_h^{r_2} B_h^{r_4} \sum_{h=1}^k \left[ (-n_1)^{\alpha_h + \beta_h} - 1 \right] A_h^{r_1} B_h^{r_3} - \sum_{h=1}^k A_h^{r_1} B_h^{r_3} \\ &\quad \sum_{h=1}^k \left[ (-n_1)^{\alpha'_h + \beta'_h} - 1 \right] A_h^{r_2} B_h^{r_4} + Q_{r_1 r_3} Q_{r_2 r_4} - Q_{(r_1 + r_2)(r_3 + r_4)}. \end{aligned} \quad (3.38)$$

Finally, combining the various equations we get the formula:

$$\begin{aligned} {}_2M_{pab} &= n^{-1} \bar{p}_{2a2b} - \bar{p}_{ab}^2 + 2(n)^{-2(a+b+1)} (a!)^2 (b!)^2 \sum_{i,h=1}^n \sum_{\alpha_h, \alpha'_h, \beta_h, \beta'_h=0}^{a,b} \sum_{k=1}^t S \\ &\quad \left\{ \sum_{h,h'=1}^k (-n_1)^{\alpha_h + \beta_h + \alpha'_h + \beta'_h} + n_k \sum_{h=1}^k [(-n_1)^{\alpha_h + \beta_h} + (-n_1)^{\alpha'_h + \beta'_h}] + C_2^{n_k} \right\} \\ &\quad \frac{\Pi P^{ih} (\alpha_h + \alpha'_h) (\beta_h + \beta'_h)}{\Pi \alpha_h! \Pi \beta_h! \Pi \alpha'_h! \Pi \beta'_h!} + 2(n)^{-2(a+b+1)} (a!)^2 (b!)^2 \sum_{j,h=1}^n \sum_{r_1, r_2, r_3, r_4=0}^{a,b} \frac{(-n)^{Sr_1} Sr_i}{\Pi r_1!} \\ &\quad \sum_{\alpha_h, \alpha'_h, \beta_h, \beta'_h=0}^{\alpha, \beta, \gamma, \delta} \sum_{k=1}^t S \left\{ \sum_{h,h'=1}^k [(-n_1)^{\alpha_h + \beta_h + \alpha'_h + \beta'_h} - 1] A_h^{r_1} B_h^{r_3} A_h^{r_2} B_h^{r_4} \right. \\ &\quad \left. - \sum_{h=1}^k [(-n_1)^{\alpha_h + \beta_h} - 1] A_h^{r_1} B_h^{r_3} \sum_{h=1}^k A_h^{r_2} B_h^{r_4} - \sum_{h=1}^k [(-n_1)^{\alpha'_h + \beta'_h} - 1] \right. \\ &\quad \left. A_h^{r_2} B_h^{r_4} \sum_{h=1}^k A_h^{r_1} B_h^{r_3} + Q_{r_2 r_4} \sum_{h=1}^k [(-n_1)^{\alpha_h + \beta_h} - 1] A_h^{r_1} B_h^{r_3} \right\} \end{aligned}$$



$$\begin{aligned}
 & + Q_{r_1 r_3} \sum_{h=1}^k [(-n_1)^{\alpha'_h + \beta'_h} - 1] A_h^{r_2} B_h^{r_4} + Q_{r_1 r_3} Q_{r_2 r_4} - Q_{(r_1 + r_2)(r_3 + r_4)} \Big\} \\
 & \frac{\Pi P_{(\alpha + \alpha'_h)(\beta_h + \beta'_h)}^{j_h}}{\Pi \alpha_h! \Pi \beta_h! \Pi \alpha'_h! \Pi \beta'_h!} \cdot
 \end{aligned} \tag{12}$$

In case the  $n$  populations are identical the second part of (12) must vanish, and in the first part the summations

$$\sum_{j_h=1}^n \prod_{h=1}^k P_{(\alpha_h + \alpha'_h)(\beta_h + \beta'_h)}^{j_h} = \frac{k! C_k^n \Pi P_{(\alpha_h + \alpha'_h)(\beta_h + \beta'_h)}}{l_1! l_2! \dots l_c!},$$

where  $l_1, l_2, \dots, l_c$  are the number of repetitions of the pairs of integers  $(\alpha_1 + \alpha'_1)(\beta_1 + \beta'_1), \dots, (\alpha_k + \alpha'_k)(\beta_k + \beta'_k)$ , respectively.

We then have the following

Corollary: The mathematical expectation of the variance,  ${}_2m_{pab}$ , of the product moment,  $p_{ab}$ , in samples of  $n$  from a single infinite population is given by

$$\begin{aligned}
 {}_2M_{pab} & = \bar{p}_{2a2b} - \bar{p}_{ab}^2 + 2(n)^{-2(a+b+1)} (a!)^2 (b!)^2 \sum_{\alpha_h, \alpha'_h, \beta_h, \beta'_h=0}^{a, a, b, b} S \\
 & \sum_{k=1}^t \frac{k! C_k^n}{l_1! l_2! \dots l_c!} \left\{ \sum_{h, h'=1}^k (-n_1)^{\alpha_h + \beta_h + \alpha'_h + \beta'_h} + n_k \sum_{h=1}^k [(-n_1)^{\alpha_h + \beta_h} \right. \\
 & \left. + (-n_1)^{\alpha'_h + \beta'_h}] + C_2^{nk} \right\} \frac{\prod_{h=1}^k P_{(\alpha_h + \alpha'_h)(\beta_h + \beta'_h)}}{\Pi \alpha_h! \Pi \beta_h! \Pi \alpha'_h! \Pi \beta'_h!} \cdot
 \end{aligned} \tag{12'}$$

4. **The Formula for  ${}_2M_{p_{21}}$ .** Formula (12) can by no means be used mechanically. It does, however, summarize to a great extent the details in finding  ${}_2M_{p_{ab}}$  for any given values  $a, b$ . Formulae for  ${}_2M_{p_{21}}, {}_2M_{p_{31}}$  have been obtained, but the one for  ${}_2M_{p_{31}}$  is too long to be included in the paper, especially since with a little work it can be easily derived by applying (12). The one for  ${}_2M_{p_{21}}$  is given immediately below.

$$\begin{aligned}
 {}_2M_{p_{21}} & = n^{-6} \{ n_1^2 n_2^2 S[P_{42}^i - (P_{21}^i)^2] + n_2^2 S[P_{40}^i P_{02}^i + 4(P_{30}^i P_{12}^i - n_2 P_{31}^i P_{11}^i)] \\
 & - 2n_2^2 n_3 S P_{22}^i P_{20}^i + (n_2^2 + 2) S(P_{20}^i P_{20}^i P_{02}^i + 8P_{20}^i P_{11}^i P_{11}^i) + 6S P_{20}^i P_{20}^i P_{02}^i \} \\
 & + 2n^{-5} \{ n_1 n_2^2 S(P_{41}^i B_i + 2P_{32}^i A_i - P_{30}^i P_{11}^i B_i - 2P_{21}^i P_{11}^i A_i) \\
 & - 4n_2 n_4 S(P_{30}^i B_i P_{11}^i + P_{12}^i A_i P_{20}^i) - 2n_2 S[n_3 P_{21}^i B_i P_{20}^i + 2(2n_2 - 3) P_{21}^i A_i P_{11}^i] \\
 & + 6n S P_{21}^i P_{11}^i A_i + 4n_2 S(P_{30}^i P_{11}^i B_i + P_{30}^i P_{02}^i A_i + P_{30}^i A_i P_{02}^i + 2P_{21}^i P_{20}^i B_i \\
 & + P_{12}^i P_{20}^i A_i) \} + n^{-4} \{ n_2^2 S[P_{40}^i B_i^2 - (P_{20}^i B_i)^2] + 4S P_{20}^i P_{20}^i (B_i + B_i)^2 \\
 & + 3(n_2^2 + n_2) S P_{22}^i A_i^2 + 4S P_{20}^i P_{02}^i (A_i + A_i)^2 - 2n_4 S[P_{20}^i P_{02}^i A_i^2
 \end{aligned}$$

$$\begin{aligned}
& + 2P_{11}^i P_{11}^j (A_i + A_j)^2] + 16SP_{11}^i A_i P_{11}^j A_j - 4n_2^2 S(P_{11}^i A_i)^2 \\
& + 4(2n_2^2 + n_2)SP_{31}^i A_i B_i - 4n_4 SP_{11}^i A_i B_i P_{20}^j - 8n_3 SP_{11} P_{20}^j A_j B_j \\
& + 8S(P_{11}^i B_i P_{20}^j A_j + P_{11}^i A_i P_{20}^j B_j) - 4n_2^2 SP_{11}^i A_i P_{20}^j B_i \\
& - 2n_1 n_2 n^{-1} S(Q_{20} P_{22}^i + 2Q_{11} P_{31}^i) + 2n_2 n^{-1} S[6Q_{11} P_{11}^i P_{20}^j \\
& + Q_{20}(P_{20}^i P_{02}^j + 4P_{11}^i P_{11}^j)] + 2n^{-4} \{nn_2 S(2P_{30}^i A_i B_i^2 + 2P_{12}^i A_i^3 + 5P_{21}^i A_i^2 B_i) \\
& - n_1 S[Q_{20}(P_{21}^i B_i + P_{12}^i A_i) + 2Q_{11} 8P_{30}^i B_i + 2P_{21}^i A_i]\} + n^{-4} \{n^2 S[P_{02}^i A_i^4 \\
& + 4(P_{20}^i A_i^2 B_i^2 + P_{11}^i A_i^2 B_i)] - 2n S[(Q_{20} A_i B_i + Q_{11} A_i^2) P_{11}^i + Q_{11} P_{20}^i A_i B_i \\
& + Q_{20} P_{02}^i A_i^2] + S[Q_{20}^2 P_{02}^i + 4Q_{20}(Q_{11} P_{11}^i + Q_{11}^2 P_{20}^i)]\}. \quad (13)^{12}
\end{aligned}$$

#### CHAPTER IV. The Mathematical Expectation of the Third Moment of $p_{11}$

1. **The Mathematical Expectation of  ${}_3m_{p_{11}}$ .** Following the notation of the last chapter we shall denote the third moment of  $p_{11}$  about its mean by  ${}_3m_{p_{11}}$  and the mathematical expectation of  ${}_3m_{p_{11}}$  by  ${}_3M_{p_{11}}$ . We have then by definition.

$${}_3m_{p_{11}} = \{n^{-1}S(x_i - x)(y_i - y) - \bar{p}_{11}\}^3,$$

and by a well known formula we have:

$${}_3M_{p_{11}} = \overline{p_{11}^3} - 3\bar{p}_{11}M_{p_{11}} - \bar{p}_{11}^3. \quad (4.11)$$

The last two terms of (4.11) are given by (1) and (11). To evaluate  $\overline{p_{11}^3}$  we write:

$$\begin{aligned}
\overline{p_{11}^3} &= E\{n^{-1}S(x_i - x)(y_i - y)\}^3 = n^{-3}SE(x_i - x)^3(y_i - y)^3 \\
&+ 3n^{-3}SE(x_i - x)^2(y_i - y)^2(x_j - x)(y_j - y) \\
&+ 6n^{-3}SE(x_i - x)(y_i - y)(x_j - x)(y_j - y)(x_k - x)(y_k - y).
\end{aligned}$$

The first term is simply  $n^{-2}\bar{p}_{33}$  which is given by (10). The evaluation of the second term is not essentially different from the evaluation of the left hand side of (3.22), and since all details have been given there we shall omit them here.

To evaluate the last expression let us write:

$$\begin{aligned}
& SE(x_i - x)(y_i - y)(x_j - x)(y_j - y)(x_k - x)(y_k - y) \\
&= SE[(U_i + A_i)(V_i + B_i)(U_j + A_j)(V_j + B_j)(U_k + A_k)(V_k + B_k)] \\
&= SE(U_i V_i U_j V_j U_k V_k) + SE(U_i V_i U_j V_j U_k B_k) + \dots + SE(A_i B_i A_j B_j A_k B_k). \quad (4.12)
\end{aligned}$$

<sup>12</sup> In case the  $n$  populations are identical this reduces to one of Pepper's formulae, *Biometrika*, Vol. XXI, p. 238, Cor. 1.

As there is a great deal of similarity among the various terms of the right hand side of (4.12), it will not be necessary to go into the details of the expansion of every one of them. We shall, therefore, indicate the details for the expansion of only two of them—one symmetrical and one non-symmetrical; and as the first two terms are of that type we shall use these for the purpose of illustration.

Using the  $u, v$  notation we have

$$SE(U_i V_i U_j V_j U_k V_k) = n^{-6} SE[(n_1 u_i - \dots)(n_1 v_i - \dots)(n_1 u_j - \dots) \\ (n_1 v_j - \dots)(n_1 u_k - \dots)(n_1 v_k - \dots)].$$

The maximum number of subscripts appearing in any term evidently being 3, we can write without any loss in generality:

$$\begin{aligned}
SE[(n_1 u_i - \dots) \dots (n_1 v_k - \dots)] &= E[(n_1 u_1 - \dots)(n_1 v_1 - \dots)(n_1 u_2 - \dots) \\
&n_1 v_2 - \dots)(n_1 u_3 - \dots)(n_1 v_3 - \dots)] + E\{(n_1 u_1 - \dots)(n_1 v_1 - \dots)[(n_1 u_2 - \dots) \\
&(n_1 v_2 - \dots) + (n_1 u_3 - \dots)(n_1 v_3 - \dots)] + (n_1 u_2 - \dots)(n_1 v_2 - \dots) \\
&(n_1 u_3 - \dots)(n_1 v_3 - \dots)\} S(n_1 u_i - \dots)(n_1 v_i - \dots) + E\{(n_1 u_1 - \dots) \\
&(n_1 v_1 - \dots) + (n_1 u_2 - \dots)(n_1 v_2 - \dots) + n_1 u_3 - \dots)(n_1 v_3 - \dots)\} \\
&S(n_1 u_i - \dots)(n_1 v_i - \dots)(n_1 u_j - \dots)(n_1 v_j - \dots) + SE\{(n_1 u_i - \dots) \dots \\
&(n_1 v_k - \dots)\} .
\end{aligned} \tag{4.13}$$

The coefficients of the various terms arising in this expansion can now be found quite easily. For example, the coefficient of  $P_{33}^1$ , which is, of course, the same as the coefficient of  $P_{33}^i$ , is easily found to be

$$n_1^2 + n_3(2n_1^2 + 1) + \frac{n_3n_4(n_1^2 + 2)}{2} + \frac{n_3n_4n_5}{6} = \frac{nn_1n_2(3n_1 - 2)}{6}.$$

To evaluate the summation  $SE(U_i V_i U_j V_j U_k B_k) = n^{-5} SE[(n_1 u_i - \dots)(n_1 v_i - \dots)(n_1 u_j - \dots)(n_1 v_j - \dots)(n_1 v_k - \dots) B_k]$ , we break it up into partial summations as follows:

$$\begin{aligned}
& SE[(n_1u_i - \dots)(n_1v_i - \dots)(n_1u_j - \dots)(n_1v_j - \dots)(n_1u_k - \dots)B_k] \\
& = E\{(n_1u_1 - \dots)(n_1v_1 - \dots)[(n_1u_2 - \dots)(n_1v_2 - \dots)(n_1u_3 - \dots)B_3 \\
& + (n_1u_2 - \dots)B_2(n_1u_3 - \dots)(n_1v_3 - \dots)] + (n_1u_1 - \dots)B_1(n_1u_2 - \dots) \\
& (n_1v_2 - \dots)(n_1u_3 - \dots)(n_1v_3 - \dots)\} + E\{(n_1u_1 - \dots)(n_1v_1 - \dots) \\
& [(n_1u_2 - \dots)(n_1v_2 - \dots) + (n_1u_3 - \dots)(n_1v_3 - \dots)] + (n_1u_2 - \dots) \\
& (n_1v_2 - \dots)(n_1u_3 - \dots)(n_1v_3 - \dots)\} S(n_1u_j - \dots)B_j + E\{(n_1u_1 - \dots) \\
& (n_1v_1 - \dots)[(n_1u_2 - \dots)B_2 + (n_1u_3 - \dots)B_3] + (n_1u_2 - \dots)(n_1v_2 - \dots)
\end{aligned}$$

$$\begin{aligned}
& [(n_1 u_1 - \dots) B_1 + (n_1 u_3 - \dots) B_3] + (n_1 u_3 - \dots)(n_1 v_3 - \dots) \\
& [(n_1 u_1 - \dots) B_1 + (n_1 v_2 - \dots) B_2] \} S(n_1 u_j - \dots)(n_1 v_j - \dots) \\
& + E\{(n_1 u_1 - \dots)(n_1 v_1 - \dots) + (n_1 u_2 - \dots)(n_1 v_2 - \dots) + (n_1 u_3 - \dots) \\
& (n_1 v_3 - \dots)\} S(n_1 u_i - \dots)(n_1 v_i - \dots)(n_1 u_j - \dots) B_j + E\{(n_1 u_1 - \dots) B_1 \\
& + (n_2 u_2 - \dots) B_2 + (n_1 u_3 - \dots) B_3\} S(n_1 u_i - \dots)(n_1 v_i - \dots) \\
& (n_1 u_j - \dots)(n_1 v_j - \dots) + ES(n_1 u_i - \dots)(n_1 v_i - \dots)(n_1 u_j - \dots) \\
& (n_1 v_j - \dots)(n_1 u_k - \dots) B_k.
\end{aligned} \tag{4.14}$$

The expansion of (4.14) is not as difficult as it appears for only two subscripts can appear in any term: the explicit appearance of the subscript 3 is due to the fact that we are dealing with a triple summation. We, consequently, do not need to expand those parentheses in which  $B$  appears.

We shall now, without any further details, state the final result, which is:

$$\begin{aligned}
{}_3M_{p11} = n^{-6} \{ & S[n_1^3 P_{33}^i - P_{30}^i P_{03}^j + 3n_1(P_{31}^i P_{02}^j + P_{20}^i P_{13}^j) + 3n_1(n_1^2 + 2)P_{22}^i P_{11}^j \\
& - 3(2n_1^2 + 1)P_{21}^i P_{12}^j + 3n_3 P_{11}^i P_{02}^j P_{20}^k + 6(n_1^3 + 3n_1 - 2)P_{11}^i P_{11}^j P_{11}^k] \\
& - 3n_1 S P_{11}^i [S(n_1^2 P_{22}^i + P_{20}^i P_{02}^j - n_1^2 (P_{11}^i)^2 + 2P_{11}^i P_{11}^j)] - n_1^3 (S P_{11}^i)^3 \} \\
& + 3n^{-5} \{ S[n_1^2 (P_{32}^i B_i + P_{23}^i A_i) + 2\alpha (P_{21}^i P_{11}^j B_i + P_{12}^i P_{11}^j A_i) \\
& - 2n_1 (P_{21}^i P_{11}^j B_i + P_{12}^i P_{11}^j A_i) - 2n_1 (P_{11}^i P_{21}^j B_i + P_{11}^i P_{12}^j A_i) \\
& + (P_{12}^i P_{20}^j B_i + P_{21}^i P_{02}^j A_i) - 2n_1 (P_{12}^i P_{20}^j B_j + P_{21}^i P_{02}^j A_j) \\
& + (P_{30}^i P_{02}^j B_j + P_{03}^i P_{20}^j A_j)] \} + 3n^{-4} \{ S[n_1 (P_{31}^i B_i^2 + P_{13}^i A_i^2) \\
& + n_2 (P_{20}^i P_{11}^j B_i^2 + P_{02}^i P_{11}^j A_i^2) - (P_{11}^i P_{20}^j B_i^2 + P_{02}^i P_{11}^j A_i^2) \\
& - 2(P_{20}^i B_i P_{11}^j B_j + P_{02}^i A_i P_{11}^j A_j) + 2n_1 P_{22}^i A_i B_i - 2P_{20}^i B_i P_{02}^j A_j \\
& - 2P_{11}^i A_i P_{11}^j B_j + 2n_2 P_{11}^i P_{11}^j A_i B_i - 2(P_{11}^i)^2 A_i B_i] \} + n^{-3} \{ S[(P_{30}^i B_i^3 + P_{03}^i A_i^3) \\
& + 3(P_{21}^i A_i B_i^2 + P_{12}^i A_i^2 B_i)] \}.
\end{aligned} \tag{14}^{13}$$

Where  $\alpha = n_1^2 + n_1 + 1$ .

This formula is shorter and simpler than the formula for  ${}_2M_{p21}$ , although they are of the same order. This is due to the symmetry of  ${}_3M_{p11}$ .

## CHAPTER V. Product Moments of Trivariate and Quadrivariate Populations

**1. Some additional definitions and notation.** In this chapter we shall indicate briefly how the method of the previous chapters may be extended to populations

<sup>13</sup> Cf. *Biometrika*, Vol. XXI, p. 253, formula (19).

of more than two variables. We shall do this by deriving some of the simpler formulae, corresponding to those of Chapter II, for trivariate and quadrivariate populations.

The notation will be slightly changed in that we shall symbolize the new variables by priming the symbols for the variables used in the previous chapters. Thus, we shall indicate the  $k^{\text{th}}$  trivariate population by  $(X_k, Y_k, X'_k)$  and the  $k^{\text{th}}$  quadrivariate population by  $(X_k, Y_k, X'_k, Y'_k)$ , and samples from such populations by  $(x_k, y_k, x'_k)$  and  $(x_k, y_k, x'_k, y'_k)$  respectively.

We shall denote by  $P_{ijk}^m$  the product moment of the  $m^{\text{th}}$  population of order  $i$  in  $X$ ,  $j$  in  $Y$ , and  $k$  in  $X'$ , and by  $P_{ijkl}^m$  the similar product moment for a quadrivariate population. These are defined by the following equations:

$$P_{ijk}^m = E(X_m - a_m)^i (Y_m - b_m)^j (X'_m - c_m)^k, \quad (5.11)$$

$$P_{ijkl}^m = E(X_m - a_m)^i (Y_m - b_m)^j (X'_m - c_m)^k (Y'_m - d_m)^l \quad (5.12)$$

where  $a_m, b_m$ , etc. are defined as in Chapter I part 2.

The sample product moments corresponding to  $P_{ijk}^m, P_{ijkl}^m$  will be denoted by  $p_{ijk}$  and  $p_{ijkl}$  respectively. They are defined by:

$$p_{ijk} = n^{-1} \sum_{m=1}^n (x_m - \bar{x})(y_m - \bar{y})(x'_m - \bar{x}')^k, \quad (5.13)$$

$$p_{ijkl} = n^{-1} \sum_{m=1}^n (x_m - \bar{x})(y_m - \bar{y})(x'_m - \bar{x}')^k (y'_m - \bar{y}')^l. \quad (5.14)$$

Finally we shall designate  $E(p_{ijk})$  and  $E(p_{ijkl})$  by  $\bar{p}_{ijk}$  and  $\bar{p}_{ijkl}$  respectively.

**2. The Mathematical Expectation of  $p_{111}$  and  $p_{211}$ .** By definition we have

$$\bar{p}_{111} = E[n^{-1} \sum_{i=1}^n S(x_i - \bar{x})(y_i - \bar{y})(x'_i - \bar{x}')]. \quad (5.21)$$

Applying the transformations (1.17) this equation becomes

$$\begin{aligned} np_{111} &= E[S(U_i + A_i)(V_i + B_i)(U'_i + C_i)] = SE(U_i V_i U'_i) + SE(U_i V_i C_i) \\ &+ SE(U_i U'_i B_i) + SE(V_i U'_i A_i) + \text{vanishing terms} + SE(A_i B_i C_i). \end{aligned} \quad (5.22)$$

Since  $EA_i B_i C_i = A_i B_i C_i$ ,  $SE(A_i B_i C_i) = SA_i B_i C_i$ . Following the previous notation we shall put  $SA_i B_i C_i = Q_{111}$ .

When the expression  $SE(U_i V_i U'_i)$  is expanded, no other non-vanishing terms except those of the form  $E(u_i v_i u'_i) = P_{111}^i$  can appear. The coefficient of this term will evidently be the same as that of  $P_{21}^i$  in (2.23), namely:  $n^{-2} n_1 n_2$ . Whence:

$$SE(U_i V_i U'_i) = n^{-2} n_1 n_2 SP_{11}^i.$$

The three terms following the first of (5.22) are by (2.24) equal to

$$n^{-1} n_2 S(P_{110}^i C_i + P_{101}^i B_i + P_{011}^i A_i).$$